



# Classical Statistics and Statistical Learning in Imaging Neuroscience

Danilo Bzdok

## ► To cite this version:

Danilo Bzdok. Classical Statistics and Statistical Learning in Imaging Neuroscience: Two Statistical Cultures in Neuroimaging. *Frontiers in Human Neuroscience*, 2017. hal-01583175

**HAL Id: hal-01583175**

**<https://hal.science/hal-01583175>**

Submitted on 6 Sep 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Classical Statistics and Statistical Learning in Imaging Neuroscience

Danilo Bzdok<sup>1,2,3</sup>

1 Department of Psychiatry, Psychotherapy and Psychosomatics, Medical Faculty, RWTH Aachen, Germany

2 JARA, Jülich-Aachen Research Alliance, Translational Brain Medicine, Aachen, Germany

3 Parietal team, INRIA, Neurospin, bat 145, CEA Saclay, 91191 Gif-sur-Yvette, France

Running title: Two Statistical Cultures in Neuroimaging

Corresponding author: danilo[DOT]bzdok[AT]rwth-aachen[DOT]de

**Abstract:** Brain-imaging research has predominantly generated insight by means of classical statistics, including regression-type analyses and null-hypothesis testing using  $t$ -test and ANOVA. Throughout recent years, statistical learning methods enjoy increasing popularity especially for applications in rich and complex data, including cross-validated out-of-sample prediction using pattern classification and sparsity-inducing regression. This concept paper discusses the implications of inferential justifications and algorithmic methodologies in common data analysis scenarios in neuroimaging. It is retraced how classical statistics and statistical learning originated from different historical contexts, build on different theoretical foundations, make different assumptions, and evaluate different outcome metrics to permit differently nuanced conclusions. The present considerations should help reduce current confusion between model-driven classical hypothesis testing and data-driven learning algorithms for investigating the brain with imaging techniques.

**Keywords:** neuroimaging, data science, epistemology, statistical inference, machine learning, p value, Rosetta Stone

## 1 Main Text

2 "The trick to being a scientist is to be open to using a wide variety of tools."

3 Leo Breiman (2001)

## 5 Introduction

6 Among the greatest challenges humans face are cultural misunderstandings between individuals, groups,  
7 and institutions (Hall, 1976). The topic of the present paper is the culture clash between knowledge  
8 generation based on null-hypothesis testing and out-of-sample pattern generalization (Breiman, 2001;  
9 Donoho, 2015; Friedman, 1998; Shmueli, 2010). These statistical paradigms are now increasingly combined  
10 in brain-imaging studies (Kriegeskorte et al., 2009; Varoquaux and Thirion, 2014). Ensuing inter-cultural  
11 misunderstandings are unfortunate because the invention and application of new research methods has  
12 always been a driving force in the neurosciences (Greenwald, 2012; Yuste, 2015). Here the goal is to  
13 disentangle the contexts underlying *classical statistical inference* and *out-of-sample generalization* by  
14 providing a direct comparison of their historical trajectories, modeling philosophies, conceptual  
15 frameworks, and performance metrics.

16 During recent years, neuroscience has transitioned from qualitative reports of few patients with  
17 neurological brain lesions to quantitative lesion-symptom mapping on the voxel level in hundreds of  
18 patients (Gläscher et al., 2012). We have gone from manually staining and microscopically inspecting single  
19 brain slices to 3D models of neuroanatomy at micrometer scale (Amunts et al., 2013). We have also gone  
20 from experimental studies conducted by a single laboratory to automatized knowledge aggregation across  
21 thousands of previously isolated neuroimaging findings (Fox et al., 2014; Yarkoni et al., 2011). Rather than  
22 laboriously collecting in-house data published in a single paper, investigators are now routinely reanalyzing  
23 multi-modal data repositories (Derrfuss and Mar, 2009; Kandel et al., 2013; Markram, 2012; Poldrack and  
24 Gorgolewski, 2014; Van Essen et al., 2012). The detail of neuroimaging datasets is hence growing in terms  
25 of information resolution, sample size, and complexity of meta-information (Bzdok and Yeo, 2017; Eickhoff  
26 et al., 2016; Van Horn and Toga, 2014). As a consequence of the data demand of many pattern-recognition  
27 algorithms, the scope of neuroimaging analyses has expanded beyond the predominance of regression-  
28 type analyses combined with null-hypothesis testing (Fig. 1). Applications of statistical learning methods i)  
29 are more data-driven due to particularly flexible models, ii) have scaling properties compatible with high-  
30 dimensional data with myriads of input variables, and iii) follow a heuristic agenda by prioritizing useful  
31 approximations to patterns in data (Blei and Smyth, 2017; Jordan and Mitchell, 2015; LeCun et al., 2015).  
32 *Statistical learning* (Hastie et al., 2001) henceforth comprises the umbrella of "machine learning", "data  
33 mining", "pattern recognition", "knowledge discovery", "high-dimensional statistics", and bears close  
34 relation to "data science".

35 From a technical perspective, one should make a note of caution that holds across application domains  
36 such as neuroscience: While the research question often precedes the choice of statistical model, perhaps

1 no single criterion exists that alone allows for a clear-cut distinction between classical statistics and  
2 statistical learning in all cases. For decades, the two statistical cultures have evolved in partly independent  
3 sociological niches (Breiman, 2001). There is currently a scarcity of scientific papers and books that would  
4 provide an explicit account on how concepts and tools from classical statistics and statistical learning are  
5 exactly related to each other. Efron and Hastie are perhaps among the first to discuss the issue in their  
6 book "Computer-Age Statistical Inference" (2016). The authors cautiously conclude that statistical learning  
7 inventions, such as support vector machines, random-forest algorithms, and "deep" neural networks, can  
8 not be easily situated in the classical theory of 20th century statistics. They go on to say that  
9 "pessimistically or optimistically, one can consider this as a bipolar disorder of the field or as a healthy  
10 duality that is bound to improve both branches" (Efron and Hastie, 2016, p. 447). In the current absence of  
11 a commonly agreed-upon theoretical account from the technical literature, the present concept paper  
12 examines applications of classical statistics versus statistical learning in the concrete context of  
13 neuroimaging analysis questions.

14 More generally, ensuring that a statistical effect discovered in one set of data extrapolates to new  
15 observations in the brain can take different forms (Efron, 2012). As one possible definition, "the goal of  
16 statistical inference is to say what we have learned about the population  $X$  from the observed data  $x$ "  
17 (Efron and Tibshirani, 1994). In a similar spirit, a committee report to the National Academies of the USA  
18 stated (Jordan et al., 2013, p. 8): "Inference is the problem of turning data into knowledge, where  
19 knowledge often is expressed in terms of variables [...] that are not present in the data per se, but are  
20 present in models that one uses to interpret the data." According to these definitions, *statistical inference*  
21 *can be understood as encompassing not only the classical null-hypothesis testing framework but also*  
22 *Bayesian model inversion to compute posterior distributions as well as more recently emerged pattern-*  
23 *learning algorithms relying on out-of-sample generalization* (cf. Cohen, 1990; Efron, 2012; Ghahramani,  
24 2015; Gigerenzer and Murray, 1987). The important consequence for the present considerations is that  
25 classical statistics and statistical learning can give rise to different categories of inferential thinking  
26 (Chamberlin, 1890; Efron and Tibshirani, 1994; Platt, 1964) - an investigator may ask an identical  
27 neuroscientific question in different mathematical contexts.

28 For a long time, knowledge generation in psychology, neuroscience, and medicine has been dominated by  
29 classical statistics with *estimation* of linear-regression-like models and subsequent *statistical significance*  
30 *testing* whether an effect exists in the sample. In contrast, computation-intensive pattern learning methods  
31 have always had a strong focus on *prediction* in frequently extensive data with more modest concern for  
32 interpretability and the "right" underlying question (Ghahramani, 2015; Hastie et al., 2001). In many  
33 statistical learning applications, it is standard practice to quantify the ability of a predictive pattern to  
34 extrapolate to other samples, possibly in individual subjects. In a two-step procedure, a learning algorithm  
35 is fitted on a bigger amount of available data (*training data*) and the ensuing fitted model is empirically  
36 evaluated on a smaller amount of independent data (*test data*). This stands in contrast to classical

1 statistical inference where the investigator seeks to reject the null hypothesis by considering the entirety of  
2 a data sample (Wasserstein and Lazar, 2016), typically all available subjects. In this case, the desired  
3 relevance of a statistical relationship in the underlying population is ensured by formal mathematical  
4 proofs and is not commonly ascertained by explicit evaluations on new data (Breiman, 2001; Wasserstein  
5 and Lazar, 2016). As such, generating insight according to classical statistics and statistical learning serves  
6 rather distinct modeling purposes. Classical statistics and statistical learning do therefore not judge data on  
7 the same aspects of evidence (Arbabshirani et al., 2017; Breiman, 2001; Bzdok and Yeo, 2017; Shmueli,  
8 2010). The two statistical cultures perform different types of principled assessment for successful  
9 extrapolation of a statistical relationship beyond the particular observations at hand.

10 Taking an epistemological perspective helps appreciating that scientific research is rarely an entirely  
11 objective process but deeply depends on the beliefs and expectations of the investigator. A new “scientific  
12 fact” about the brain is probably not established in vacuo (Fleck et al., 1935; terms in quotes taken from  
13 source). Rather, a research “object” is recognized and accepted by the “subject” according to socially  
14 conditioned “thought styles” that are cultivated among members of “thought collectives”. A witnessed and  
15 measured neurobiological phenomenon tends to only become “true” if not at odds with the constructed  
16 “thought history” and “closed opinion system” shared by that subject. The present paper will revisit and  
17 reintegrate two such thought milieus in the context of imaging neuroscience: classical statistics (CISt) and  
18 statistical learning (StLe).

## 20 **Different histories: The origins of classical hypothesis testing and pattern-learning algorithms**

21 One of many possible ways to group statistical methods is by framing them along the lines of CISt and StLe.  
22 The incongruent historical developments of the two statistical communities are even evident from their  
23 basic terminology. Inputs to statistical models are usually called *independent variables*, *explanatory*  
24 *variables*, or *predictors* in the CISt community, but are typically called *features* collected in a *feature space*  
25 in the StLe community. The model outputs are typically called *dependent variables*, *explained variable*, or  
26 *responses* in CISt, while these are often called *target variables* in StLe. It follows a summary of characteristic  
27 events in the development of what can today be considered as CISt and StLe (Fig. 2).

28 Around 1900 the notions of *standard deviation*, *goodness of fit*, and the  $p < 0.05$  threshold emerged  
29 (Cowles and Davis, 1982). This was also the period when William S. Gosset published the *t-test* under the  
30 incognito name “Student” to quantify production quality in Guinness breweries. Motivated by concrete  
31 problems such as the interaction between potato varieties and fertilizers, Ronald A. Fisher invented the  
32 *analysis of variance* (ANOVA), *null-hypothesis testing*, promoted *p values*, and devised principles of proper  
33 experimental conduct (Fisher, 1925; Fisher, 1935; Fisher and Mackenzie, 1923). Another framework by  
34 Jerzy Neyman and Egon S. Pearson proposed the *alternative hypothesis*, which allowed for the statistical  
35 notions of *power*, *false positives* and *false negatives*, but left out the concept of *p values* (Neyman and  
36 Pearson, 1933). This was a time before electrical calculators emerged after World War II (Efron and

1 Tibshirani, 1991; Gigerenzer, 1993). Student's *t*-test and Fisher's inference framework were  
2 institutionalized by American psychology textbooks widely read in the 40s and 50s, while Neyman and  
3 Pearson's framework only became increasingly known in the 50s and 60s. Today's applied statistics  
4 textbooks have inherited a mixture of the Fisher and Neyman-Pearson approaches to statistical inference.

5 It is a topic of current debate<sup>1,2,3</sup> whether ClSt is a discipline that is separate from StLe (e.g., Bishop and  
6 Lasserre, 2007; Breiman, 2001; Chambers, 1993; Efron and Hastie, 2016; Friedman, 2001; Shalev-Shwartz  
7 and Ben-David, 2014) or if "statistics" denotes a broader methodological class that includes both ClSt and  
8 StLe tools as its members (e.g., Blei and Smyth, 2017; Cleveland, 2001; Jordan and Mitchell, 2015; Tukey,  
9 1962). StLe methods may be more often adopted by computer scientists, physicists, engineers, and others  
10 who typically have less formal statistical background and may be more frequently working in industry  
11 rather than academia. In fact, John W. Tukey foresaw many of the developments that led up to what one  
12 might today call statistical learning (Tukey, 1962). He early proposed a "peaceful collision of computing and  
13 statistics". A modern reformulation of the same idea states (Efron and Hastie, 2016): "If the  
14 inference/algorithm race is a tortoise-and-hare affair, then modern electronic computation has bred a  
15 bionic hare." Indeed, kernel methods, decision trees, nearest-neighbor algorithms, graphical models, and  
16 various other statistical tools actually emerged in the ClSt community, but largely continued to develop in  
17 the StLe community (Friedman, 2001).

18 As often cited beginnings of statistical learning approaches, the *perceptron* was an early brain-inspired  
19 computing algorithm (Rosenblatt, 1958), and Arthur Samuel created a checker board program that  
20 succeeded in beating its own creator (Samuel, 1959). Such studies towards *artificial intelligence* (AI) led to  
21 enthusiastic optimism and subsequent periods of disappointment during the so-called "AI winters" in the  
22 late 70s and around the 90s (Cox and Dean, 2014; Kurzweil, 2005; Russell and Norvig, 2002), while the  
23 increasingly available computers in the 80s encouraged a new wave of statistical algorithms (Efron and  
24 Tibshirani, 1991). Later, the use of StLe methods increased steadily in many quantitative scientific domains  
25 as they underwent an increase in data richness from classical "long data" (samples  $n > \text{variables } p$ ) to  
26 increasingly encountered "wide data" ( $n \ll p$ ) (Hastie et al., 2015; Tibshirani, 1996). The emerging field of  
27 StLe has received conceptual consolidation by the seminal book "The Elements of Statistical Learning"  
28 (Hastie et al., 2001). The coincidence of changing data properties, increasing computational power, and  
29 cheaper memory resources encouraged a still ongoing resurgence in StLe research and applications  
30 approximately since 2000 (Manyika et al., 2011; UK House of Common, 2016). For instance, over the last 15  
31 years, *sparsity* assumptions gained increasing relevance for statistical and computational tractability as well  
32 as for domain interpretability when using *supervised* and *unsupervised* learning algorithms (i.e., with and

---

<sup>1</sup> "Data Science and Statistics: different worlds?" (Panel at Royal Statistical Society UK, March 2015)  
(<https://www.youtube.com/watch?v=C1zMUjHOLr4>)

<sup>2</sup> "50 years of Data Science" (David Donoho, Tukey Centennial workshop, USA, September 2015)

<sup>3</sup> "Are ML and Statistics Complementary?" (Max Welling, 6th IMS-ISBA meeting, December 2015)

1 without target variables) in the high-dimensional " $n \ll p$ " setting (Bühlmann and Van De Geer, 2011; Hastie  
2 et al., 2015). More recently, improvements in training very "deep" (i.e., many non-linear hidden layers)  
3 neural-networks architectures (Hinton and Salakhutdinov, 2006) have much improved automatized feature  
4 selection (Bengio et al., 2013) and have exceeded human-level performance in several application domains  
5 (LeCun et al., 2015).

6 In sum, "the biggest difference between pre- and post-war statistical practice is the degree of automation"  
7 (Efron and Tibshirani, 1994) up to a point where "almost all topics in twenty-first-century statistics are now  
8 computer-dependent" (Efron and Hastie, 2016). ClSt has seen many important inventions in the first half of  
9 the 20th century, which have often developed at statistical departments of academic institutions and  
10 remain in nearly unchanged form in current textbooks of psychology and other empirical sciences. The  
11 emergence of StLe as a coherent field has mostly taken place in the second half of the 20th century as a  
12 number of disjoint developments in industry and often non-statistical departments in academia (e.g., AT&T  
13 Bell Laboratories), which lead for instance to artificial neural networks, support vector machines, and  
14 boosting algorithms (Efron and Hastie, 2016). Today, systematic education in StLe is still rare at the large  
15 majority of universities, in contrast to the many consistently offered ClSt courses (Burnham and Anderson,  
16 2014; Cleveland, 2001; Donoho, 2015; Vanderplas, 2013).

17 In neuroscience, the advent of brain-imaging techniques, including positron emission tomography (PET) and  
18 functional magnetic resonance imaging (fMRI), allowed for the in-vivo characterization of the neural  
19 correlates underlying sensory, cognitive, or affective tasks. Brain scanning enabled *quantitative* brain  
20 measurements with *many variables per observation* (analogous to the advent of high-dimensional  
21 microarrays in genetics; Efron, 2012). Since the inception of PET and fMRI, deriving topographical  
22 localization of neural activity changes was dominated by analysis approaches from ClSt, especially the  
23 general linear model (GLM; Poline and Brett, 2012; Scheffé, 1959). The classical approach to neuroimaging  
24 analysis is probably best exemplified by the statistical parametric mapping (SPM) software package that  
25 implements the GLM to provide a mass-univariate characterization of regionally specific effects.

26 As distributed information over voxels is less well captured by many ClSt approaches, including common  
27 GLM applications, StLe models were proposed early on for neuroimaging investigations. For instance,  
28 principal component analysis was used to distinguish globally distributed neural activity changes (Moeller  
29 et al., 1987) as well as to study Alzheimer's disease (Grady et al., 1990). Canonical correlation analysis was  
30 used to quantify complex relationships between task-free neural activity and schizophrenia symptoms  
31 (Friston et al., 1992). However, these first approaches to "multivariate" brain-behavior associations did not  
32 ignite a major research trend (cf. Friston et al., 2008; Worsley et al., 1997). The popularity of StLe methods  
33 in neuroimaging only peaked after being rebranded as "mind-reading", "brain decoding", and "multivariate  
34 pattern analysis" or "MVPA" that appealed by identifying ongoing thought from neural activity (Haynes and  
35 Rees, 2005; Kamitani and Tong, 2005). Up to that point, the term *prediction* had less often been used by  
36 imaging neuroscientists in the sense of out-of-sample generalization of a learning algorithm and more often

in the incompatible sense of (in-sample) linear correlation such as using Pearson's or Spearman's method (Gabrieli et al., 2015; Shmueli, 2010). While there was scarce discussion of the position of "decoding" models in formal statistical terms, growing interest was manifested in first review publications and tutorial papers on applying StLe methods to neuroimaging data (Haynes and Rees, 2006; Mur et al., 2009; Pereira et al., 2009). The interpretational gains of this new access to the neural representation of behavior and its disturbances in disease was flanked by the availability of necessary computing power and memory resources. Although challenging to realize, "deep" neural network algorithms have recently been introduced to neuroimaging research (de Brebisson and Montana, 2015; Güçlü and van Gerven, 2015; Plis et al., 2014). These computation-intensive models might help in approximating and deciphering the nature of neural processing in brain circuits (Cox and Dean, 2014; Yamins and DiCarlo, 2016). As the dimensionality and complexity of neuroimaging datasets are constantly increasing, neuroscientific investigations will be always more likely to benefit from StLe methods given their natural scaling to large-scale data analysis (Blei and Smyth, 2017; Efron, 2012; Efron and Hastie, 2016).

From a conceptual viewpoint (Fig. 3), a large majority of statistical methods can be situated somewhere on a continuum between the two poles of ClSt and StLe (Efron and Hastie, 2016; Jordan et al., 2013; p. 61). ClSt was mostly fashioned for problems with small samples that can be grasped by plausible models with a small number of parameters chosen by the investigator in an analytical fashion. StLe was mostly fashioned for problems with many variables in potentially large samples with little knowledge of the data-generating process that gets emulated by a mathematical function derived from data in a heuristic fashion. Tools from ClSt therefore typically assume that the data behave according to certain known mechanisms, whereas StLe exploits algorithmic techniques to avoid many a-priori specifications of data-generating mechanisms. Neither ClSt or StLe nor any of the other categories of statistical models can be considered generally superior. This relativism is captured by the so-called *no free lunch theorem*<sup>4</sup> (Wolpert, 1996): no single statistical strategy can consistently do better in all circumstances (cf. Gigerenzer, 2004). As a very general rule of thumb, ClSt preassumes and formally tests *a model for the data*, whereas StLe extracts and empirically evaluates *a model from the data*.

### **Case study one: Cognitive contrast analysis and decoding mental states**

Vignette: A neuroimaging investigator wants to reveal the neural correlates underlying face processing in humans. 40 healthy, right-handed adults are recruited and undergo a block design experiment run in a 3T MRI scanner with whole-brain coverage. In a passive viewing paradigm, 60 colored stimuli of unfamiliar faces are presented, which have forward head and gaze position. The control condition presents colored

---

<sup>4</sup> In the supervised setting, there is no a priori distinction between learning algorithms evaluated by out-of-sample prediction error. In the optimization setting of finite spaces, all algorithms searching an extremum perform identically when averaged across possible cost functions. (<http://www.no-free-lunch.org/>)



1 pictures of 60 different houses to the participants. In the experimental paradigm, a picture of a face or a  
2 house is presented for 2 seconds in each trial and the interval between trials within each block is randomly  
3 jittered varying from 2 to 7 seconds. The picture stimuli are presented in pseudo-randomized fashion and  
4 are counterbalanced in each passively watching participant. Despite the blocked presentation of stimuli,  
5 each experiment trial is modeled separately. The fMRI data are analyzed using a GLM as implemented in  
6 the SPM software package. Two task regressors are included in the model for the face and house conditions  
7 based on the stimulus onsets and viewing durations and using a canonical hemodynamic response function.  
8 In the GLM design matrix, the face column and house column are hence set to 1 for brain scans from the  
9 corresponding task condition and set to 0 otherwise. Separately in each brain voxel, the GLM parameters  
10 are estimated, which fits  $\beta_{\text{face}}$  and  $\beta_{\text{house}}$  regression coefficients to explain the contribution of each  
11 experimental task to the neural activity increases and decreases observed in that voxel. A  $t$ -test can then  
12 formally assess whether the fMRI signal in the current voxel is significantly more involved in viewing faces  
13 as opposed to the house control condition. ...

14 Question: What is the statistical difference between *subtracting* the neural activity from the face versus  
15 house conditions and *decoding* the neural activity during face versus house processing?

16  
17 Computing cognitive contrasts is a CISt approach that was and still is routinely performed in the *mass-*  
18 *univariate regime*: it fits a separate GLM model for each individual voxel in the brain scans and then tests  
19 for significant differences between the obtained condition coefficients (Friston et al., 1994). Instead,  
20 decoding cognitive processes from neural activity is a StLe approach that is typically performed in a  
21 *multivariate regime*: a learning algorithm is trained on a large number of voxel observations in brain scans  
22 and then the model's prediction accuracy is evaluated on sets of new brain scans. These CISt and StLe  
23 approaches to identifying the neural correlates underlying cognitive processes of interest are closely  
24 related to the notions of *encoding models* and *decoding models*, respectively (Kriegeskorte, 2011; Naselaris  
25 et al., 2011; Pedregosa et al., 2015; but see Güçlü et al, 2015).

26 Encoding models regress the brain data against a design matrix with indicators of the face versus house  
27 condition and formally test whether the difference is statistically significant. Decoding models typically aim  
28 to predict these indicators by training and empirically evaluating classification algorithms on different splits  
29 from the whole dataset. In CISt parlance, the model *explains* the neural activity, the *dependent or explained*  
30 *variable*, measured in each separate brain voxel, by the *beta coefficients* according to the experimental  
31 condition indicators in the *design matrix* columns, the *independent or explanatory variables*. That is, the  
32 GLM can be used to explain neural activity changes by a linear combination of experimental variables  
33 (Naselaris et al., 2011). Answering the same neuroscientific question with decoding models in StLe jargon,  
34 the *model weights* of a *classifier* are fitted on the *training set* of the *input data* to *predict* the *class labels*,  
35 the *target variables*, and are subsequently evaluated on the *test set* by *cross-validation* to obtain their *out-*  
36 *of-sample generalization performance*. Here, classification algorithms are used to predict entries of the

design matrix by identifying a linear or more complicated combination between the many simultaneously considered brain voxels (Pereira et al., 2009). More broadly, CISt applications in functional neuroimaging tend to estimate the presence of cognitive processes from neural activity, whereas many StLe applications estimate properties of neural activity from different cognitive tasks.

A key difference between many CISt-mediated encoding models and StLe-mediated decoding models thus pertains to the direction of statistical estimation between brain space and behavior space (Friston et al., 2008; Varoquaux and Thirion, 2014). It was noted (Friston et al., 2008) that the direction of brain-behavior association is related to the question whether the stimulus indicators in the model act as causes by representing deterministic experimental variables of an encoding model or consequences by representing probabilistic outputs of a decoding model. Such considerations also reveal the intimate relationship of CISt models to the notion of *forward inference*, while StLe methods are probably more often used for formal *reverse inference* in functional neuroimaging (Eickhoff et al., 2011; Poldrack, 2006; Varoquaux and Thirion, 2014; Yarkoni et al., 2011). On the one hand, *forward inference* relates to encoding models by testing the probability of observing activity in a brain location given knowledge of a psychological process. On the other hand, *reverse inference* relates to brain decoding to the extent that classification algorithms can learn to distinguish experimental fMRI data to belong to two psychological conditions and subsequently be used to estimate the presence of specific cognitive processes based on new neural activity observations (cf. Poldrack, 2006). Finally, establishing a brain-behavior association has been argued to be more important than the actual direction of the mapping function (Friston, 2009). This author stated that "showing that one can decode activity in the visual cortex to classify [...] a subject's percept is exactly the same as demonstrating significant visual cortex responses to perceptual changes" and, conversely, "all demonstrations of functionally specialized responses represent an implicit mindreading".

Conceptually, GLM-based encoding models follow a *representational agenda* by testing hypotheses on *regional effects of functional specialization* in the brain (where?). A *t*-test is used to compare pairs of neural activity estimates to statistically distinguish the target face and the non-target house condition (Friston et al., 1996). Essentially, this test for significant differences between the fitted beta coefficients corresponds to two stimulus indicators based on well-founded arguments from cognitive theory. This statistical approach assumes that *cognitive subtraction* is possible, that is, the regional brain responses of interest can be isolated by contrasting two sets of brain scans that are believed to differ in the cognitive facet of interest (Friston et al., 1996; Stark and Squire, 2001). For one voxel location at a time, an attempt is made to reject the null hypothesis of no difference between the averaged *neural activity level* of a target brain state and the averaged neural activity of a control brain state. It is important to appreciate that the representational agenda thus emphasizes the *relative difference* in fMRI signal during tasks and may neglect the individual neural activity information of each particular task (Logothetis et al., 2001). Note that the univariate GLM analysis can be extended to more than one output (dependent or explained) variable

1 within the CISt regime by performing a multivariate analysis of covariance (MANCOVA). This allows for tests  
 2 of more complex hypotheses but incurs multivariate normality assumptions (Kriegeskorte, 2011).  
 3 More generally, it is seldom mentioned that the GLM would not have been solvable for unique solutions in  
 4 the high-dimensional " $n \ll p$ " regime, instead of fitting one model for each voxel in the brain scans. This is  
 5 because the number of brain voxels  $p$  exceed by far the number of data samples  $n$  (i.e., leading to an  
 6 under-determined system of equations), which incapacitates many statistical estimators from CISt (cf.  
 7 Giraud, 2014; Hastie et al., 2015). Regularization by sparsity-inducing norms, such as in modern *penalized*  
 8 regression analysis using the LASSO and ElasticNet, emerged only later (Tibshirani, 1996; Zou and Hastie,  
 9 2005) as a principled StLe strategy to de-escalate the need for dimensionality reduction or preliminary  
 10 filtering of important voxels and to enable the tractability of the high-dimensional analysis setting.  
 11 Consequently, many software packages for the analysis of cognitive neuroimaging experiments have  
 12 implemented discrete voxel-wise analyses with classical inference instead of more recent StLe alternatives.  
 13 Because hypothesis testing for significant differences between beta coefficients of fitted GLMs relies on  
 14 comparing the means of neural activity measurements, the results from statistical tests are not corrupted  
 15 by the conventionally applied spatial smoothing with a Gaussian filter. On the contrary, this image  
 16 preprocessing step even helps the correction for multiple comparisons based on random fields theory (cf.  
 17 below), alleviates inter-individual neuroanatomical variability, and can thus increase sensitivity. Spatial  
 18 smoothing however discards fine-grained neural activity patterns spatially distributed across voxels that  
 19 potentially carry information associated with mental operations (cf. Haynes, 2015; Kamitani and Sawahata,  
 20 2010). Indeed, some authors believe that sensory, cognitive, and motor processes manifest themselves as  
 21 "neuronal population codes" (Averbeck et al., 2006). Relevance of such population codes in human  
 22 neuroimaging was for instance suggested by revealing subject-specific neural responses in the fusiform  
 23 gyrus to facial stimuli (Saygin et al., 2012). In applications of StLe models, the spatial smoothing step is  
 24 therefore often skipped because the "decoding" algorithms precisely exploit the locally varying structure of  
 25 the salt-and-pepper patterns in fMRI signals.  
 26 In so doing, decoding models use learning algorithms in an *informational agenda* by showing *generalization*  
 27 *of robust patterns* to new brain activity acquisitions (de-Wit et al., 2016; Kriegeskorte et al., 2006; Mur et  
 28 al., 2009). Information that is weak in one voxel but spatially distributed across voxels can be effectively  
 29 harvested in a structure-preserving fashion (Haynes, 2015; Haynes and Rees, 2006). This modeling agenda  
 30 is focused on the whole *neural activity pattern*, in contrast to the representational agenda dedicated to  
 31 separate increases or decreases in *neural activity level*. For instance, the default mode network typically  
 32 exhibits activity *decreases* at the onset of many psychological tasks with visual or other sensory stimuli,  
 33 whereas the induced activity *patterns* in that less activated network may nevertheless functionally subserve  
 34 task execution (Bzdok et al., 2016; Christoff et al., 2016). Some brain-behavior associations might only  
 35 emerge when simultaneously capturing neural activity in a group of voxels but disappear in single-voxel  
 36 approaches, such as mass-univariate GLM analyses (cf. Davatzikos, 2004). Note that, analogous to

multivariate variants of the GLM, decoding could also be replaced by classical statistical approaches. Furthermore, inference on extrapolating information patterns in the brain reduces to model comparison (Friston, 2009). During training of a classification algorithm to predict face versus house stimuli based on many brain voxels, an optimization algorithm (e.g., gradient descent or its variants) searches iteratively through the *hypothesis space* (= *function space*) of the chosen learning model. Each such hypothesis corresponds to one specific combination of model weights (i.e., a weighted contribution of individual brain measurements) that equates with one candidate mapping function from the neural activity features to the target variables indicating face and house stimulation.

Among other views, it has previously been proposed (Brodersen, 2009) that four types of neuroscientific questions become readily quantifiable through StLe applications to neuroimaging: i) *Where* is an information category neurally processed? This can extend the interpretational spectrum from increase and decrease of neural activity to the existence of complex combinations of activity variations distributed across voxels. For instance, linear classifiers could decode object categories from the ventral temporal cortex even after excluding the fusiform gyrus, which is known to be responsive to object stimuli (Haxby et al., 2001). ii) *Whether* a given information category is reflected by neural activity? This can extend the interpretational spectrum to topographically similar but neurally distinct processes that potentially underlie different cognitive facets. For instance, linear classifiers could successfully distinguish whether a subject is attending to the first or second of two simultaneously presented stimuli (Kamitani and Tong, 2005). iii) *When* is an information category generated (i.e., onset), processed (i.e., duration), and bound (i.e., alteration)? When applying classifiers to neural time series, the interpretational spectrum can be extended to the beginning, evolution, and end of distinct cognitive facets. For instance, different classifiers have been demonstrated to map the decodability time structure of mental operation sequences (King and Dehaene, 2014). iv) More controversially, *how* is an information category neurally processed? The interpretational spectrum can be extended to computational properties of the neural processes, including processing in brain regions versus brain networks or isolated versus partially shared processing facets. For instance, a classifier trained for evolutionarily conserved eye gaze movement was able to decode evolutionarily more recent mathematical calculation processes as a possible case of “neural recycling” in the human brain (Anderson, 2010; Knops et al., 2009). As an important caveat in interpreting StLe models, the particular technical properties of a chosen learning algorithm (e.g., linear versus non-linear support vector machines) can probably seldom serve as a convincing argument for reverse-engineering mechanisms of neural information processing as measured by fMRI scanning (cf. Misaki et al., 2010).

In sum, the statistical properties of ClSt and StLe methods have characteristic consequences in neuroimaging analysis and interpretation. They hence offer different access routes and complementary answers to identical neuroscientific questions.

## 1 Case study two: Small volume correction and searchlight analysis

2 Vignette: The neuroimaging experiment from case study 1 successfully identified the fusiform gyrus of the  
3 ventral visual stream to be more responsive to face stimuli than house stimuli. However, the investigator's  
4 initial hypothesis of also observing face-responsive neural activity in the ventromedial prefrontal cortex  
5 could not be confirmed in the *whole-brain* analyses. The investigator therefore wants to follow up with a  
6 *topographically focused* approach that examines differences in neural activity between the face and house  
7 conditions exclusively in the ventromedial prefrontal cortex.

8 Question: What are the statistical implications of delineating task-relevant neural responses in a spatially  
9 constrained search space rather than analyzing brain measurements of the entire brain?

10

11 A popular CIST approach to corroborate less pronounced neural activity findings is *small volume correction*.  
12 This region of interest (ROI) analysis involves application of the mass-univariate GLM approach only to the  
13 ventromedial prefrontal cortex as a preselected biological compartment, rather than considering the grey-  
14 matter voxels of the entire brain in a naïve, topographically unconstrained fashion. Small volume correction  
15 allows for significant findings in the ROI that remain sub-threshold after accounting for the tens of  
16 thousands of multiple comparisons in the whole-brain GLM analysis. Small volume correction is therefore a  
17 simple means to alleviate the multiple-comparisons problem that motivated more than two decades of still  
18 ongoing methodological developments in the neuroimaging domain (Friston, 2006; Nichols, 2012; Smith et  
19 al., 2001; Worsley et al., 1992). Whole-brain GLM results were initially reported as uncorrected findings  
20 without accounting for multiple comparisons, then with Bonferroni's family wise error (FWE) correction,  
21 later by random field theory correction using neural activity height (or clusters), followed by false discovery  
22 rate (FDR) (Genovese et al., 2002) and slowly increasing adoption of cluster-thresholding for voxel-level  
23 inference via permutation testing (Smith and Nichols, 2009). Rather than the isolated voxel, it has early  
24 been discussed that a possibly better unit of interest should be spatially neighboring voxel groups (see here  
25 for an overview: Chumbley and Friston, 2009). The setting of high regional correlation of neural activity was  
26 successfully addressed by random field theory that provide inferences not about individual voxels but  
27 topological features in the underlying (spatially continuous) effects. This topological inference is used to  
28 identify clusters of relevant neural activity changes from their peak, size or mass (Worsley et al., 1992).  
29 Importantly, the spatial dependencies of voxel observations were not incorporated into the GLM  
30 estimation step, but instead taken into account during the subsequent model inference step to alleviate the  
31 multiple-comparisons problem.

32 A cousin of small volume correction in the StLe world would be to apply classification algorithms to a subset  
33 of voxels to be considered as input to the model (i.e., *feature selection*). In particular, *searchlight analysis* is  
34 an increasingly popular learning technique that can identify *locally constrained multivariate patterns* in  
35 neural activity (Friman et al., 2001; Kriegeskorte et al., 2006). For each voxel in the ventromedial prefrontal

1 cortex, the brain measurements of the immediate neighborhood are first collected (e.g., radius of 10mm  
2 voxels). In each such searchlight, a classification algorithm, for instance linear support vector machines, is  
3 then trained on one part of the brain scans (*training set*) and subsequently applied to determine the  
4 prediction accuracy in the remaining, unseen brain scans (*test set*). In this StLe approach, the excess of  
5 brain voxels is handled by performing pattern recognition analysis in only dozens of locally adjacent voxel  
6 neighborhoods at a time. Finally, the mean classification accuracy of face versus house stimuli across all  
7 permutations over the brain data is mapped to the center of each considered sphere. The searchlight is  
8 then moved through the ROI until each seed voxel had once been the center voxel of the searchlight. This  
9 yields a voxel-wise classification map of accuracy estimates for the entire ventromedial prefrontal cortex.  
10 Consistent with the informational agenda (cf. above), searchlight analysis quantifies the extent to which  
11 (local) neural activity *patterns* can *predict* the difference between the house and face conditions. It  
12 contrasts small volume correction that determines whether one experimental condition exhibited a  
13 significant neural activity *increase* or *decrease* relative to a particular other experimental condition,  
14 consistent with the representational agenda.

15 When considering high-dimensional brain scans through the ClSt lens, the statistical challenge resides in  
16 solving the *multiple-comparisons problem* (Nichols, 2012; Nichols and Hayasaka, 2003). From the StLe  
17 stance, however, it is the *curse of dimensionality* and *overfitting* that statistical analyses need to tackle  
18 (Domingos, 2012; Friston et al., 2008). Many neuroimaging analyses based on ClSt methods can be viewed  
19 as testing a particular hypothesis (i.e., the null hypothesis) repeatedly in a large number of separate voxels.  
20 In contrast, testing whether learning algorithm extrapolate to new brain data can be viewed as searching  
21 through thousands of different hypotheses in a single process (i.e., walking through the hypothesis space;  
22 cf. above) (Shalev-Shwartz and Ben-David, 2014).

23 As common brain scans offers measurements of >100,000 brain locations, a mass-univariate GLM analysis  
24 typically entails the same statistical test to be applied >100,000 times. The more often the investigator tests  
25 a hypothesis of relevance for a brain location, the more locations will be falsely detected as relevant (false  
26 positive, Type I error), especially in the noisy neuroimaging data. All dimensions in the brain data (i.e., voxel  
27 variables) are implicitly treated as equally important and no neighborhoods of most expected variation are  
28 statistically exploited (Hastie et al., 2001). Hence, the absence of restrictions on observable structure in the  
29 set of data variables during the statistical modeling of neuroimaging data takes a heavy toll at the final  
30 inference step. This is where *random field theory* comes to the rescue. As noted above, this form of  
31 topological inference dispenses with the problem of inferring which voxels are significant and tries to  
32 identify significant topological features in the underlying distributed responses. By definition, topological  
33 features like maxima are sparse events and can be thought of as a form of dimensionality reduction - not in  
34 data space but in the statistical characterization of where neural responses occur.

35 This is contrasted by the high-dimensional StLe regime, where the initial model family chosen by the  
36 investigator determines the complexity restrictions to all data dimensions (i.e., all voxels, not single voxels)

that are imposed explicitly or implicitly by the model structure. Model choice predisposes existing but unknown low-dimensional neighborhoods in the full voxel space to achieve the prediction task. Here, the toll is taken at the beginning of the investigation because there are so many different alternative model choices that would impose a different set of complexity constraints to the high-dimensional measurements in the brain. For instance, signals from "brain regions" are likely to be well approximated by models that impose discrete, locally constant compartments on the data (e.g., k-means or spatially constrained Ward clustering). Instead, tuning model choice to signals from macroscopical "brain networks" should impose overlapping, locally continuous data compartments (e.g., independent component analysis or sparse principal component analysis) (Bzdok et al., 2017; Bzdok and Yeo, 2017; Yeo et al., 2014).

Exploiting such *effective dimensions* in the neuroimaging data (i.e., coherent brain-behavior associations involving many distributed brain voxels) is a rare opportunity to simultaneously reduce the *model bias* and *model variance*, despite their typical inverse relationship (Hastie et al., 2001). Model bias relates to prediction failures incurred because the learning algorithm can systematically not represent certain parts of the underlying relationship between brain scans and experimental conditions (formally, the deviation between the target function and the average function space of the model). Model variance relates to prediction failures incurred by noise in the estimation of the optimal brain-behavior association (formally, the difference between the best-choice input-output relation and the average function space of the model). A model that is too simple to capture a brain-behavior association probably underfits due to high bias. Yet, an overly complex model probably overfits due to high variance. Generally, high-variance approaches are better at *approximating* the "true" brain-behavior relation (i.e., in-sample model estimation), while high-bias approaches have a higher chance of *generalizing* the identified pattern to new observations (i.e., out-of-sample model evaluation). The bias-variance tradeoff can be useful in explaining why successful applications of statistical models largely rely on i) the amount of available data, ii) the typically not known amount of noise in the data, and iii) the unknown complexity of the target function in nature (Abu-Mostafa et al., 2012).

Learning algorithms that overcome the curse of dimensionality - extracting coherent patterns from all brain voxels at once - typically incorporate an implicit bias for anisotropic neighborhoods in the data (Bach, 2014; Bzdok et al., 2015; Hastie et al., 2001). Put differently, prediction models successful in the high-dimensional setting have an in-built specialization to representing types of functions that are compatible with the structure to be uncovered in the brain data. Knowledge embodied in a learning algorithm suited to a particular application domain can better calibrate the sweet spot between underfitting and overfitting. When applying a model without any complexity restrictions to high-dimensional data generalization becomes difficult to impossible because all directions in the data (i.e., individual brain voxels) are treated equally with isotropic structure. At the root of the problem, all data samples look virtually identical to the learning algorithm in high-dimensional data scenarios (Bellman, 1961). The learning algorithm will not be able to see through the idiosyncracies in the data, will tend to overfit, and thus be unlikely to generalize to

new observations. Such considerations provide insight into why the multiple-comparisons problem is more often an issue in encoding studies, while overfitting is more closely related to decoding studies (Friston et al., 2008). The juxtaposition of ClSt and StLe views offers insights into why restricting neural data analysis to a ROI with fewer voxels, rather than the whole brain, simultaneously alleviates both the multiple-comparisons problem (ClSt) and the curse of dimensionality (StLe).

As a practical summary, drawing classical inference in neuroimaging data has largely been performed by considering each voxel independently and by massive simultaneous testing of a same null hypothesis in all observed voxels. This has incurred a multiple-comparisons problem difficult enough that common approaches may still be prone to incorrect results (Efron, 2012). In contrast, aiming for generalization of a pattern in high-dimensional neuroimaging data to new observations in the brain incurs the equally challenging curse of dimensionality. Successfully accounting for the high number of input dimensions will probably depend on learning models that impose neurobiologically justified bias and keeping the variance under control by dimensionality reduction and regularization techniques.

More broadly, asking at what point new neurobiological knowledge is arising during ClSt and StLe investigations relies on largely distinct theoretical frameworks that revolve around *null-hypothesis testing* and *statistical learning theory* (Fig. 4). Both ClSt and StLe methods share the common goal of demonstrating relevance of a given effect in the data beyond the sample brain scans at hand. However, the attempt to show successful extrapolation of a statistical relationship at the general population is embedded in different mathematical contexts. Knowledge generation in ClSt and StLe is hence rooted in different notions of statistical inference.

ClSt laid down its most important inferential framework in the Popperian spirit of critical empiricism (Popper, 1935/2005): scientific progress is to be made by continuous replacement of current hypotheses by ever more pertinent hypotheses using *falsification*. The rationale behind hypothesis falsification is that one counterexample can reject a theory by *deductive reasoning*, while any quantity of evidence can not confirm a given theory by inductive reasoning (Goodman, 1999). The investigator verbalizes two mutually exclusive hypotheses by domain-informed judgment. The *alternative hypothesis* should be conceived as the outcome intended by the investigator and to contradict the state of the art of the research topic. The *null hypothesis* represents the devil's advocate argument that the investigator wants to reject (i.e., falsify) and it should automatically deduce from the newly articulated alternative hypothesis. A conventional 5%-threshold (i.e., equating with roughly two standard deviations) guards against rejection due to the idiosyncrasies of the sample that are not representative of the general population. If the data have a probability of  $\leq 5\%$  given the null hypothesis ( $P(\text{result} | H_0)$ ), it is evaluated to be significant. Such a *test for statistical significance* indicates a difference between two means with a 5% chance of being a false positive finding. If the null hypothesis can not be rejected (which depends on power), then the test yields no conclusive result, rather than a null result (Schmidt, 1996). In this way, classical hypothesis testing continuously replaces currently embraced hypotheses explaining a phenomenon in nature by better hypotheses with more empirical



support in a Darwinian selection process. Finally, Fisher, Neyman, and Pearson intended hypothesis testing as a marker for further investigation, rather than an off-the-shelf decision-making instrument (Cohen, 1994; Nuzzo, 2014).

In StLe instead, answers to how neurobiological conclusions can be drawn from a dataset at hand are provided by the *Vapnik-Chervonenkis dimensions* (VC dimensions) from *statistical learning theory* (Vapnik, 1989, 1996). The VC dimensions of a pattern-learning algorithm quantify the probability at which the distinction between the neural correlates underlying the face versus house conditions can be captured and used for correct predictions in new, possibly later acquired brain scans from the same cognitive experiment (i.e., *out-of-sample generalization*). Such statistical approaches implement the *inductive* strategy to learn general principles (i.e., the neural signature associated with given cognitive processes) from a series of exemplary brain measurements, which contrasts the *deductive* strategy of rejecting a certain null hypothesis based on counterexamples (cf. Bengio, 2014; Lake et al., 2015; Tenenbaum et al., 2011). The VC dimensions measure how complicated the examined relationship between brain scans and experimental conditions could become - in other words, the richness of the representation which can be instantiated by the used model, the complexity capacity of its *hypothesis space*, the “wiggleness” of the decision boundary used to distinguish examples from several classes, or, more intuitively, the “currency” of learnability. VC dimensions are derived from the maximal number of different brain scans that can be correctly detected to belong to either the house condition or the face condition by a given model. The VC dimensions thus provide a theoretical guideline for the largest set of brain scan examples fed into a learning algorithm such that this model is able to guarantee zero classification errors.

As one of the most important results from statistical learning theory, in any intelligent learning system, the opportunity to derive abstract patterns in the world by reducing the discrepancy between prediction error from training data (in-sample estimate) and prediction error from independent test data (out-of-sample estimate) decreases with the higher model capacity and increases with the number of available training observations (Vapnik, 1996; Vapnik and Kotz, 1982). In brain imaging, a learning algorithm is hence theoretically backed up to successfully predict outcomes in future brain scans with high probability if the chosen model ignores structure that is overly complicated, such as higher-order non-linearities between many brain voxels, and if the model is provided with a sufficient number of training brain scans. Hence, VC dimensions provide explanations why increasing the number of considered brain voxels as input features (i.e., entailing increased number of model parameters) or using a more sophisticated prediction model, requires more training data for successful generalization. Notably, the VC dimensions (analogous to null-hypothesis testing) are unrelated to the *target function*, as the “true” mechanisms underlying the studied phenomenon in nature. Nevertheless, the VC dimensions provide justification that a certain learning model can be used to approximate that target function by fitting a model to a collection of input-output pairs. In short, VC dimensions is among the best frameworks to derive theoretical errors bounds for predictive models (Abu-Mostafa et al., 2012).

1 Further, some common invalidations of the ClSt and StLe statistical frameworks are conceptually related.  
2 An often-raised concern in neuroimaging studies performing classical inference is *double dipping* or *circular*  
3 *analysis* (Kriegeskorte et al., 2009). This occurs when, for instance, first correlating a behavioral measure  
4 with brain activity and then using the identified subset of brain voxels for a second correlation analysis with  
5 that same behavioral measurement (Lieberman et al., 2009; Vul et al., 2008). In this scenario, voxels are  
6 submitted to two statistical tests with the same goal in a nested, non-independent fashion<sup>5</sup> (Freedman,  
7 1983). This corrupts the *validity of the null hypothesis* on which the reported test results conditionally  
8 depend. Importantly, this case of repeating a same statistical estimation with iteratively pruned data  
9 selections (on the training data split) is a valid routine in the StLe framework, such as in recursive feature  
10 extraction (Guyon et al., 2002; Hanson and Halchenko, 2008). However, double-dipping or circular analysis  
11 in ClSt applications to neuroimaging data have an analog in StLe analyses aiming at out-of-sample  
12 generalization: *data-snooping* or *peeking* (Abu-Mostafa et al., 2012; Fithian et al., 2014; Pereira et al.,  
13 2009). This occurs, for instance, when performing simple (e.g., mean-centering) or more involved (e.g., k-  
14 means clustering) target-variable-dependent or -independent preprocessing on the entire dataset if it  
15 should be applied separately to the training sets and test sets. Data-snooping can lead to overly optimistic  
16 cross-validation estimates and a trained learning algorithm that fails on fresh data drawn from the same  
17 distribution (Abu-Mostafa et al., 2012). Rather than a corrupted null hypothesis, it is the *error bounds of the*  
18 *VC dimensions that are loosened* and, ultimately, invalidated because information from the concealed test  
19 set influences model selection on the training set.

20 As a conceptual summary, statistical inference in ClSt is drawn by using the *entire data* at hand to *formally*  
21 *test for theoretically guaranteed* extrapolation of an effect to the general population. In stark contrast,  
22 inferential conclusions in StLe are typically drawn by fitting a model on a *larger part of the data* at hand  
23 (i.e., in-sample model selection) and *empirically testing* for successful extrapolation to an independent,  
24 smaller part of the data (i.e., out-of-sample model evaluation). As such, ClSt has a focus on *in-sample*  
25 *estimates* and *explained-variance* metrics that measure some form of goodness of fit, while StLe has a focus  
26 on *out-of-sample estimates* and *prediction accuracy*.

27

### 28 **Case study three: Significant group differences and predicting the group of participants**

29 Vignette: After isolating the neural correlates underlying face processing, the neuroimaging investigator  
30 wants to examine their relevance in psychiatric disease. In addition to the 40 healthy participants, 40  
31 patients diagnosed with schizophrenia are recruited and administered the same experimental paradigm  
32 and set of face and house pictures. In this clinical fMRI study on group differences, the investigator wants  
33 to explore possible imaging-derived markers that index deficits in social-affective processing in patients  
34 carrying a diagnosis of schizophrenia.

---

<sup>5</sup> "If you torture the data enough, nature will always confess." (Ronald Coase)

1 Question: Can metrics of statistical relevance from ClSt and StLe be combined to corroborate a given  
2 candidate biomarker?

3  
4 Many investigators in imaging neuroscience share a background in psychology, biology, or medicine, which  
5 includes training in traditional "textbook" statistics. Many neuroscientists have thus adopted a natural habit  
6 of assessing the quality of statistical relationships by means of p values, effect sizes, confidence intervals,  
7 and statistical power. These are ubiquitously taught and used at many universities, although they are not  
8 the only coherent set of statistical diagnostics (Fig. 5). These outcome metrics from ClSt may for instance be  
9 less familiar to some scientists with a background in computer science, physics, engineering, or philosophy.  
10 As an equally legitimate and internally coherent, yet less widely known diagnostic toolkit from the StLe  
11 community, prediction accuracy, precision, recall, confusion matrices, F1 score, and learning curves can  
12 also be used to measure the relevance of statistical relationships (Abu-Mostafa et al., 2012; Yarkoni and  
13 Westfall, 2016).

14 On a general basis, applications of ClSt and StLe methods may not judge findings on identical grounds  
15 (Breiman, 2001; Lo et al., 2015; Shmueli, 2010). There is an often-overlooked misconception that models  
16 with high explanatory performance do necessarily exhibit high predictive performance (Lo et al., 2015; Wu  
17 et al., 2009; Yarkoni and Westfall, 2016). For instance, brain voxels in ventral visual stream found to well  
18 *explain* the difference between face processing in healthy and schizophrenic participants based on an  
19 ANOVA may not in all cases be the best brain features to train a support vector machine to *predict* this  
20 group effect in new participants. An important outcome measure in ClSt is the quantified *significance*  
21 associated with a statistical relationship between few variables given a pre-specified model. ClSt tends to  
22 *test for a particular structure* in the brain data based on *analytical guarantees*, in form of as mathematical  
23 convergence theorems about approximating the population properties with increasing sample size. The  
24 outcome measure for StLe is the quantified *generalization of patterns* between many variables or, more  
25 generally, the robustness of special structure in the data (Hastie et al., 2001). In the neuroimaging  
26 literature, reports of statistical outcomes have previously been noted to confuse diagnostic measures from  
27 classical statistics and statistical learning (Friston, 2012).

28 For neuroscientists adopting a ClSt culture computing p values takes a central position. The *p value* denotes  
29 the probability of observing a result at least as extreme as a test statistic, assuming the null hypothesis is  
30 true. Results are considered significant when it is equal or below a pre-specified value, like  $p = 0.05$   
31 (Anderson et al., 2000). Under the condition of sufficiently high power (cf. below), it quantifies the strength  
32 of evidence against the null hypothesis as a continuous function (Rosnow and Rosenthal, 1989).  
33 Counterintuitively, it is not an immediate judgment on the alternative hypothesis  $H_1$  preferred by the  
34 investigator (Anderson et al., 2000; Cohen, 1994). P values do also not qualify the possibility of replication.  
35 It is another important caveat that a finding in the brain becomes more statistically significant (i.e., lower p  
36 value) with increasing sample size (Berkson, 1938; Miller et al., 2016).

1 The essentially binary p value (i.e., significant versus not significant) is therefore often complemented by  
2 continuous *effect size* measures for the importance of rejecting  $H_0$ . The effect size allows the identification  
3 of marginal effects that pass the statistical significance threshold but are not practically relevant in the real  
4 world. The p value is a deductive *inferential* measure, whereas the effect size is a *descriptive* measure that  
5 follows neither inductive nor deductive reasoning. The (normalized) effect size can be viewed as the  
6 strength of a statistical relationship - how much  $H_0$  deviates from  $H_1$ , or the likely presence of an effect in  
7 the general population (Chow, 1998; Ferguson, 2009; Kelley and Preacher, 2012). This diagnostic measure is  
8 often unit-free, sample-size independent, and typically standardized. As a property of the actual statistical  
9 test, the effect size can be essential to report for biological understanding, but has different names and  
10 takes various forms, such as *rho* in Pearson correlation, *eta*<sup>2</sup> in explained variances, and *Cohen's d* in  
11 differences between group averages.

12 Additionally, the certainty of a *point estimate* (i.e., the outcome is a value) can be expressed by an *interval*  
13 *estimate* (i.e., the outcome is a value range) using *confidence intervals* (Casella and Berger, 2002). These  
14 variability diagnostics indicate a range of values between which the true value will fall a given proportion of  
15 the time (Cumming, 2009; Estes, 1997; Nickerson, 2000). Typically, a 95% confidence interval is spanned  
16 around the population mean in 19 out of 20 cases across all observed samples. The tighter the confidence  
17 interval, the smaller the variance of the point estimate of the population parameter in each drawn sample.  
18 The estimation of confidence intervals is influenced by sample size and population variability. Confidence  
19 intervals may be asymmetrical (ignored by Gaussianity assumptions; Efron, 2012), can be reported for  
20 different statistics and with different percentage borders. Notably, they can be used as a viable surrogate  
21 for formal tests of statistical significance in many scenarios (Cumming, 2009).

22 Some confidence intervals can be computed in various data scenarios and statistical regimes, whereas the  
23 *power* may be especially meaningful within the culture of classical hypothesis testing (Cohen, 1977, 1992;  
24 Oakes, 1986). To estimate power the investigator needs to specify the true effect size and variance under  
25  $H_1$ . The CISt-minded investigator can then estimate the probability for rejecting null hypotheses that should  
26 be rejected, at the given threshold alpha and given that  $H_1$  is true. A high power thus ensures that  
27 statistically significant and non-significant tests indeed reflect a property of the population (Chow, 1998).  
28 Intuitively, a small confidence interval around a relevant effect suggests high statistical power. False  
29 negatives (i.e., Type II errors, beta error) become less likely with higher power (= 1 - beta error) (cf.  
30 Ioannidis, 2005). Concretely, an underpowered investigation means that the investigator is less likely to be  
31 able to distinguish between  $H_0$  and  $H_1$  at the specified significance threshold alpha. Power calculations  
32 depend on several factors, including significance threshold alpha, the effect size in the population, variation  
33 in the population, sample size n, and experimental design (Cohen, 1992).

34 While neuroimaging studies based on classical statistical inference ubiquitously report p values and  
35 confidence intervals, there have however been few reports of effect size in the neuroimaging literature  
36 (Kriegeskorte et al., 2010). Effect sizes are however necessary to compute power estimates. This explains

1 the even rarer occurrence of power calculations in the neuroimaging literature (but see Poldrack et al.,  
2 2017; Yarkoni and Braver, 2010). Given the importance of p values *and* effect sizes, the goal of computing  
3 both these useful statistics, such as for group differences in the neural processing of face stimuli, can be  
4 achieved based on two independent samples of these experimental data (especially if some selection  
5 process has been used). One sample would be used to perform statistical inference on the neural activity  
6 change yielding a p value and one sample to obtain unbiased effect sizes. Further, it has been previously  
7 emphasized (Friston, 2012) that p values and effect sizes reflect in-sample estimates in a retrospective  
8 inference regime (CISt). These metrics find an analogue in out-of-sample estimates issued from cross-  
9 validation in a prospective prediction regime (StLe). In-sample effect sizes are typically an *optimistic*  
10 estimate of the "true" effect size (inflated by high significance thresholds), whereas out-of-sample effect  
11 sizes are *unbiased* estimates of the "true" effect size.

12 In the high-dimensional scenario, the StLe-minded investigator analyzing "wide" neuroimaging data in our  
13 case, computing and judging statistical significance by p values can become challenging (Bühlmann and Van  
14 De Geer, 2011; Efron, 2012; James et al., 2013). Instead, *classification accuracy* on fresh data is probably  
15 the most often-reported performance metric in neuroimaging studies using learning algorithms. The  
16 *classification accuracy* is a simple summary statistic that captures the fraction of correct prediction  
17 instances among all performed applications of a fitted model. Basing interpretation on accuracy alone can  
18 be an insufficient diagnostic because it is frequently influenced by the number of samples, the local  
19 characteristics of hemodynamic responses, efficiency of experimental design, data folding into train and  
20 test sets, and differences in the feature number  $p$  (Haynes, 2015). A potentially under-exploited data-driven  
21 tool in this context is *bootstrapping*. The archetypical example of computer-intensive statistical method  
22 enables population-level inference of unknown distributions largely independent of model complexity by  
23 repeated random draws from the neuroimaging data sample at hand (Efron, 1979; Efron and Tibshirani,  
24 1994). This opportunity to equip various point estimates by an interval estimate of certainty (e.g., the  
25 possibly asymmetrical interval for the "true" accuracy of a classifier) is unfortunately seldom embraced in  
26 neuroimaging today (but see Bellec et al., 2010; Pernet et al., 2011; Vogelstein et al., 2014). Besides  
27 providing confidence intervals, bootstrapping can also perform non-parametric null hypothesis testing. This  
28 may be one of few examples of a direct connection between CISt and StLe methodology. Alternatively,  
29 *binomial tests* have been used to obtain a p value estimate of statistical significance from accuracies and  
30 other performance scores (Brodersen et al., 2013; Hanke et al., 2015; Pereira et al., 2009) in the binary  
31 classification setting. It has frequently been employed to reject the null hypothesis that two categories  
32 occur equally often. There are however increasing concerns about the validity of this approach if statistical  
33 independence between the performance estimates (e.g., prediction accuracies from each cross-validation  
34 fold) is in question (Jamalabadi et al., 2016; Noirhomme et al., 2014; Pereira and Botvinick, 2011). Yet  
35 another option to derive p values from classification performances of two groups is *label permutation*  
36 based on non-parametric resampling procedures (Golland and Fischl, 2003; Nichols and Holmes, 2002). This

1 algorithmic significance-testing tool can serve to reject the null hypothesis that the neuroimaging data do  
2 not contain relevant information about the group labels in many complex data analysis settings.

3 The neuroscientist who adopted a StLe culture is in the habit of corroborating prediction accuracies using  
4 *cross-validation*: the de facto standard to obtain an unbiased estimate of a model's capacity to generalize  
5 beyond the brain scans at hand (Bishop, 2006; Hastie et al., 2001). *Model assessment* is done by training on  
6 a bigger subset of the available data (i.e., *training set* for *in-sample performance*) and subsequent  
7 application of the trained model to the smaller remaining part of data (i.e., *test set* for *out-of-sample*  
8 *performance*), both assumed to be drawn from the same distribution. Cross-validation typically permutes  
9 over the sample in data splits until the class label (i.e., healthy versus schizophrenic) of each data point has  
10 been predicted once. The pairs of model-predicted label and the corresponding true label for each data  
11 point (i.e., brain scan) in the dataset can then be submitted to the quality measures (Powers, 2011),  
12 including *prediction accuracy* (inversely related to *prediction error*), *precision*, *recall*, and *F1 score*. Accuracy  
13 and the other performance metrics are often computed separately on the training set and the test set.  
14 Additionally, the measures from training and testing can be expressed by their inverse (e.g., *training error*  
15 as *in-sample error* and *test error* as *out-of-sample error*) because the positive and negative cases are  
16 interchangeable.

17

18 **Table 1**

Notion	Formula
<i>Specificity</i>	true negative / (true negative + false positive)
<i>Sensitivity / Recall</i>	true positive / (true positive + false negative)
<i>Precision</i>	true positive / (true positive + false positive)

19

20 The classification accuracy can be further decomposed into group-wise metrics based on the so-called  
21 *confusion matrix*, the juxtaposition of the true and predicted group memberships. The *precision* measures  
22 (Table 1) how many of the labels predicted from brain scans are correct, that is, how many participants  
23 predicted to belong to a certain class really belong to that class. Put differently, among the participants  
24 predicted to suffer from schizophrenia, how many have really been diagnosed with that disease? On the  
25 other hand, the *recall* measures how many labels are correctly predicted, that is, how many members of a  
26 class were predicted to really belong to that class. Hence, among the participants known to be affected by  
27 schizophrenia, how many were actually detected as such? Precision can be viewed as a measure of  
28 "exactness" and recall as a measure of "completeness" (Powers, 2011).

29 Neither accuracy, precision, or recall allow injecting subjective importance into the evaluation process of  
30 the learning algorithm. This disadvantage is addressed by the  $F_{\text{beta}}$  score: a weighted combination of the  
31 precision and recall prediction scores. Concretely, the  $F_1$  score would equally weigh precision and recall of

class predictions, while the  $F_{0.5}$  score puts more emphasis on precision and the  $F_2$  score more on recall. Moreover, applications of recall, precision, and  $F_{\text{beta}}$  scores have been noted to ignore the true negative cases as well as to be highly susceptible to estimator bias (Powers, 2011). Needless to say, no single outcome metric can be equally optimal in all contexts.

Extending from the setting of healthy-diseased classification to the *multi-class setting* (e.g., comparing healthy, schizophrenic, bipolar, and autistic participants) injects ambiguity into the interpretation of accuracy scores. Rather than reporting mere better-than-chance findings in StLe analyses, it becomes more important to evaluate the  $F_1$ , precision and recall scores for each class to be predicted in the brain scans (e.g., Brodersen et al., 2011b; Schwartz et al., 2013). It is important to appreciate that the sensitivity/specificity metrics, perhaps more frequently reported in ClSt communities, and the precision/recall metrics, probably more frequently reported in StLe communities, tell slightly different stories about identical neuroscientific findings. In fact, sensitivity equates with recall. Specificity does however not equate with precision. Further, a ClSt view on the StLe metrics would be that maximum precision corresponds to absent Type I errors (i.e., no false positives), whereas maximum recall corresponds to absent false negatives (i.e., no Type II errors). Again, Type I and II errors are related to the entirety of data points in a ClSt regime and prediction is only evaluated on a test data split of the sample in an StLe regime. Moreover, many empirical sciences usually aggregate results in *ROC* (receiver operating characteristic) curves plotting sensitivity against specificity scores, whereas other scientific domains tend to report analogous yet different *recall-precision curves* instead (Altman and Bland, 1994; Davis and Goadrich, 2006; Demšar, 2006).

Finally, StLe-minded investigators use *learning curves* (Abu-Mostafa et al., 2012; Murphy, 2012) as an important diagnostic tool for empirical estimates of the *sample complexity*, that is, the achieved model fit and prediction accuracy as a function of the available sample size  $n$ . For increasingly bigger subsets of the training set, a classification algorithm is trained on that current share of the training set and then evaluated for accuracy on the always-same test set. Across subset instances, simple models display relatively high in-sample error because they can not approximate the target function very well (underfitting) but exhibit good generalization to unseen data with relatively low out-of-sample error. Yet, complex models display relatively low in-sample error because they adapt too well to the data (overfitting) with difficulty to extrapolate to newly sampled data with high out-of-sample error. Put differently, a big gap between high in-sample and low out-of-sample performance is typically observed for high-variance models, such as artificial neural network algorithms or random forests. These performance metrics from different data splits often converge for high-bias models, such as linear support vector machines and logistic regression.

In sum, the ClSt and StLe communities rely on diagnostic metrics that are largely incongruent and may therefore not lend themselves for direct comparison in all practical analysis settings.

#### **Case study four: Out-of-sample generalization and subsequent classical inference**

Vignette: The investigator is interested in potential differences in brain volume that are associated with an individual's age (*continuous target variable*). A LASSO (*often considered as StLe arsenal*) is computed on the voxel-based morphometry data from the brain's grey matter of the 1200-subject HCP release (Human Connectome Project; Van Essen et al., 2012). This L1-penalized residual-sum-of-squares regression performs variable selection (i.e., *effectively eliminates coefficients by setting them to zero*) on all grey-matter voxels' volume information in a *high-dimensional regime* (i.e., no mass-univariate analysis). Assessing *generalization* performance of different sparse models using five-fold cross-validation yields the non-zero coefficients for few brain voxels whose volumetric information is most *predictive* of an individual's age.

Question: How can the investigator perform *classical inference* to know which of the grey-matter voxels selected to be predictive for biological age are *statistically significant*?

This is an important concern because most statistical methods currently applied to large datasets perform some explicit or implicit form of variable selection (Hastie et al., 2015; Jenatton et al., 2011; Jordan et al., 2013). There are even many different forms of preliminary selection of variables before performing significance tests on them. First, LASSO is a widely used estimator in engineering, compressive sensing, various "omics" branches and other sciences, where it is often applied without an additional significance test. Beyond neuroscience, generalization-approved statistical learning models are routinely solving a diverse set of real-world challenges. This includes algorithmic trading in financial markets, fraud detection in credit card transactions, real-time speech translation, SPAM filtering for e-mails, face recognition in digital cameras, and piloting self-driving cars (Jordan and Mitchell, 2015; LeCun et al., 2015). In all these examples, statistical learning algorithms successfully generalize to unseen, later acquired data and thus tackle the problem heuristically without classical significance test on specific variables or for overall model performance.

Second, the LASSO has been introduced as an elegant solution to the combinatorial problem of what subset of grey-matter voxels is sufficient for predicting an individual's age by *automatic variable selection*. Computing voxel-wise p values would recast this high-dimensional pattern-learning setting (i.e., considering all brain voxels at once) into a mass-univariate hypothesis-testing problem (i.e., considering one voxel after the other) where relevance would be computed independently for each voxel and correction for multiple comparisons would become necessary. Yet, recasting into the mass-univariate setting would ignore the sophisticated selection process that led to the predictive model with a reduced number of variables (Wu et al., 2009). Put differently, the variable selection procedure is itself a stochastic process that is however not accounted for by the theoretical guarantees of classical inference for statistical significance (Berk et al., 2013). Put in yet another way, data-driven model selection is corrupting the null hypothesis of classical statistical inference because the sampling distribution of the parameter estimates is altered. The important consequence is that naive classical inference expects a non-adaptive model chosen before data acquisition



1 and can therefore not be readily used along LASSO in particular or arbitrary selection procedures in  
2 general<sup>6</sup>.

3 Third, the portrayed conflict between more exploratory model selection by cross-validation (StLe) and more  
4 confirmatory classical inference (CISt) is currently at the frontier of statistical development (Loftus, 2015;  
5 Taylor and Tibshirani, 2015). New methods for so-called *post-selection inference* (or *selective inference*)  
6 allow computing p values for a set of features that have previously been chosen to be meaningful  
7 predictors by some criterion, one example being sparsity-incuding prediction algorithms such as LASSO.  
8 According to the theory of CISt, the statistical model is to be chosen before visiting the data. Classical  
9 statistical tests and confidence intervals therefore become invalidated and the p values become  
10 optimistically biased (Berk et al., 2013). Consequently, the association between a predictor and the target  
11 variable must be even stronger to certify on the same level of significance. Selective inference for modern  
12 adaptive regression thus replaces loose *naïve p values* by more rigorous *selection-adjusted p values*. As an  
13 ordinary null hypothesis can hardly be adopted in this adaptive testing setting, conceptual extension is also  
14 prompted on the level of CISt theory itself (Hastie et al., 2015). For instance, closed-form solutions to  
15 adjusted classical inference after variable selection already exist for principal component analysis (Choi et  
16 al., 2014) and forward stepwise regression (Taylor et al., 2014). Moreover, a simple alternative to formally  
17 account for preceeding model selection is *data splitting* (Cox, 1975; Fithian et al., 2014; Wasserman and  
18 Roeder, 2009), which is frequent practice in genetics (e.g., Sladek et al., 2007). In this procedure, the  
19 variable selection procedure is computed on one data split and p values are computed on the remaining  
20 second data split. However, such data splitting is not always possible and will incur power losses.

## 22 **Case study five: Classical inference and subsequent out-of-sample generalization**

23 Vignette: The investigator is interested in potential brain structure differences that are associated with an  
24 individual's gender (*categorical target variable*) in the voxel-based morphometry data of the 1200-subject  
25 HCP release (Human Connectome Project; Van Essen et al., 2012). First, the >100,000 voxels per brain scan  
26 are reduced to the most important 10,000 voxels to lower the computational cost and facilitate estimation  
27 of a prediction model. To this end, ANOVA (*univariate test for statistical significance belonging to CISt*) is  
28 initially used to obtain a ranking of the most relevant 10,000 features from the grey matter. This selects the  
29 10,000 out of the original >100,000 voxel variables with highest variance explaining volume differences  
30 between males and females (i.e., *the gender information associated with each brain scan is used in the*  
31 *univariate test*). Second, support vector machine classification ("*multivariate*" *pattern-learning algorithm*  
32 *belonging to StLe*) is performed by cross-validation on a feature space with the 10,000 preselected grey-  
33 matter measurements to predict the gender from each subject's brain scan.

---

<sup>6</sup> "Once applied only to the selected few, the interpretation of the usual measures of uncertainty do not remain intact directly, unless properly adjusted." (Yoav Benjamini)

1 Question: Is an analysis pipeline with *univariate classical inference* and subsequent *high-dimensional*  
2 *prediction* valid if both steps rely on gender as the target variables?

3  
4 The implications of feature engineering procedures applied before training a learning algorithm is a  
5 frequent concern and can require subtle answers (Guyon and Elisseeff, 2003; Hanke et al., 2015;  
6 Kriegeskorte et al., 2009; Lemm et al., 2011). In most applications of predictive models the large majority of  
7 brain voxels will be uninformative (Brodersen et al., 2011a). The described scenario of *dimensionality*  
8 *reduction* by feature selection to focus prediction is clearly allowed under the condition that the ANOVA is  
9 not computed on the entire data sample. Rather, the initial identification of voxels explaining most variance  
10 between the male and female individuals should be computed only on the training set in each cross-  
11 validation fold. In the training set and test set of each fold the same identified candidate voxels are then  
12 regrouped into a feature space that is fed into the support vector machine algorithm. This ensures an  
13 identical feature space for model training and model testing but its construction only depends on structural  
14 brain scans from the training set. Generally, voxel preprocessing performed before model training is  
15 authorized if the feature space construction is not influenced by properties of the concealed test set. In the  
16 present scenario, the Vapnik-Chervonenkis bounds of the cross-validation estimator are therefore not  
17 loosened or invalidated if class labels have been exploited for feature selection or depending on whether  
18 the feature selection procedure is univariate or multivariate (Abu-Mostafa et al., 2012; Shalev-Shwartz and  
19 Ben-David, 2014). Put differently, the cross-validation procedure simply evaluates the entire prediction  
20 process including the automatized and potentially nested dimensionality reduction approaches. In sum, in  
21 an StLe regime, using class information during feature preprocessing for a cross-validated supervised  
22 estimator is not an instance of *data-snooping* (or *peeking*) if done exclusively on the training set (Abu-  
23 Mostafa et al., 2012).

24 At the core of this explanation is the goal of cross-validation to yield *out-of-sample estimates*. In stark  
25 contrast, remember that null-hypothesis testing yields *in-sample estimates* as it needs all available data  
26 points to take its decision. Using the class labels for a variable selection step just before null-hypothesis  
27 testing on a same data sample would invalidate the null hypothesis (Kriegeskorte et al., 2010; Kriegeskorte  
28 et al., 2009). Consequently, in a ClSt regime, using class information to select variables before null-  
29 hypothesis testing will incur an instance of *double-dipping* (or *circular analysis*). This also occurs when, for  
30 instance, first correlating a behavioral measure with brain activity and then using the identified subset of  
31 brain voxels for a second correlation analysis with that same behavioral measurement (Lieberman et al.,  
32 2009; Vul et al., 2008). In this scenario, voxels are submitted to two statistical tests with the same goal in a  
33 nested, non-independent fashion (Freedman, 1983). This corrupts the *validity of the null hypothesis* on  
34 which the reported test results conditionally depend.

35 Regarding interpretation of the results, the classifier will miss some brain voxels that only carry relevant  
36 information when considered in voxel ensembles. This is because the ANOVA filter has kept voxels that are

1 independently relevant (Brodersen et al., 2011a). Univariate feature selection in high-dimensional brain  
2 scans may therefore systematically encourage model selection (i.e., each weight combination equates with  
3 a model hypothesis from the classifier's function space) that is not tuned to neurobiological  
4 meaningfulness. Concretely, in the discussed scenario the classifier learns *complex patterns between voxels*  
5 *that were previously chosen to be individually important*. This may considerably weaken the interpretability  
6 and conclusions on "whole-brain multivariate patterns". Remember also that variables that have a  
7 *statistically significant association* with a target variable do not necessarily have good *generalization*  
8 *performance*, and vice versa (Bzdok and Yeo, 2017; Lo et al., 2015; Shmueli, 2010). On the upside, it is  
9 frequently observed that the combination of whole-brain univariate feature selection and linear  
10 classification is among the best approaches if the primary goal is maximizing *prediction performance* as  
11 opposed to maximizing *interpretability*.

12 Finally, it is interesting to consider that ANOVA-mediated feature selection to a subset of  $p < 500$  voxel  
13 variables would reduce the "wide" neuroimaging data (" $n \ll p$ " setting) down to "long" neuroimaging data  
14 with fewer features than observations (" $n > p$ " setting) given the  $n = 500$  subjects (Wainwright, 2014). This  
15 allows recasting the StLe regime into a ClSt regime in order to fit a GLM and perform classical statistical  
16 tests instead of training a predictive classification algorithm (Brodersen et al., 2011a).

17

## 18 **Case study six: Structure discovery by clustering algorithms**

19 Vignette: Each functionally specialized region in the human brain probably has a unique set of long-range  
20 connections (Passingham et al., 2002). This notion has prompted connectivity-based parcellation methods  
21 in neuroimaging that segregate a ROI (can be locally circumscribed or brain global; Eickhoff et al., 2015) into  
22 distinct cortical modules (Behrens et al., 2003). The whole-brain connectivity for each ROI voxel is  
23 computed and the voxel-wise connectional fingerprints are submitted to a clustering algorithm (i.e.,  
24 *individual brain voxels in the ROI are the elements to group; the connectivity strength values are the*  
25 *features of each element for similarity assessment*). The investigator wants to apply connectivity-based  
26 parcellation to the fusiform gyrus to segregate this ROI into cortical modules that exhibit similar  
27 connectivity patterns and are, thus potentially, functionally distinct. That is, voxels within the same cluster  
28 in the ROI will have more similar connectivity properties than voxels from different ROI clusters.

29 Question: Is it possible to decide whether the obtained brain *clusters* are *statistically significant*?

30

31 In essence, the aim of connectivity-guided brain parcellation is to find useful, simplified structure by  
32 imposing circumscribed compartments on brain topography (Frackowiak and Markram, 2015; Smith et al.,  
33 2013; Yeo et al., 2011). This is typically achieved by using k-means, hierarchical, Ward, or spectral clustering  
34 algorithms (Eickhoff et al., 2015; Thirion et al., 2014). Putting on the ClSt hat, a ROI clustering result would  
35 be deemed statistically significant if the obtained data are incompatible with the null hypothesis that the  
36 investigator seeks to reject (Everitt, 1979; Halkidi et al., 2001). Choosing a test statistic for clustering

1 solutions to obtain p values is difficult (Vogelstein et al., 2014) because of the need to find a meaningful  
 2 null hypothesis to test against (Jain et al., 1999). Put differently, for classical inference based on statistical  
 3 hypothesis testing one may need to pick an arbitrary null hypothesis to falsify. It follows that neither the  
 4 ClSt notions of effect size and power do seem to apply in the case of brain parcellation (also a frequent  
 5 question by paper reviewers). Instead of classical inference to formally *test* for a particular structure in the  
 6 clustering results, the investigator actually needs to resort to exploratory approaches that discover and  
 7 assess structure in the neuroimaging data (Efron and Tibshirani, 1991; Hastie et al., 2001; Tukey, 1962).  
 8 Although statistical methods span a continuum between the two poles of ClSt and StLe, finding a clustering  
 9 model with the highest fit in the sense of explaining the regional connectivity differences at hand is perhaps  
 10 more naturally situated in the StLe community.

11 Putting on the StLe hat, the investigator realizes that the problem of brain parcellation constitutes an  
 12 *unsupervised* learning setting without any target variable  $y$  to predict (e.g., cognitive tasks, the age or  
 13 gender of the participants). The learning problem does therefore not consist in estimating a supervised  
 14 predictive model  $y = f(X)$ , but to estimate an unsupervised descriptive model for the connectivity data  $X$   
 15 themselves. Solving such unsupervised estimation problems is generally recognized to be ill-posed because  
 16 it is generally unclear what the best way is to quantify how well relevant structure has been captured and  
 17 what notion of “relevance” is most pertinent (Bishop, 2006; Ghahramani, 2004; Hastie et al., 2001; Shalev-  
 18 Shwartz and Ben-David, 2014). In clustering analysis, there are many possible transformations, projections,  
 19 and compressions of  $X$  but there is no unique criterion of optimality that clearly suggests itself. On the one  
 20 hand, the “true” *shape of clusters* is unknown for most real-world clustering problems, including brain  
 21 parcellation studies. On the other hand, finding an “optimal” *number of clusters* represents an unresolved  
 22 issue (*cluster validity problem*) in statistics in general and in brain neuroimaging in particular (Handl et al.,  
 23 2005; Jain et al., 1999). In other words, “the clustering problem is inherently ill posed, in the sense that  
 24 there is no single criterion that measures how well a clustering of data corresponds to the real world”  
 25 (Goodfellow et al., 2016). Evaluating the adequacy of clustering results is therefore conventionally  
 26 addressed by applying different *cluster validity criteria* (Eickhoff et al., 2015; Thirion et al., 2014). These  
 27 heuristic metrics are useful and necessary because clustering algorithms will always find some subregions  
 28 in the investigator's ROI, that is, find relevant structure with respect to the particular optimization objective  
 29 of the clustering algorithm whether such structure truly exists in nature or not. The various clustering  
 30 validity criteria, possibly based on information theory, topology, or consistency (Eickhoff et al., 2015),  
 31 typically encourage cluster solutions with low within-cluster and high between-cluster differences  
 32 according to a certain notion of optimality. Given that the notions of optimality are not coherent with each  
 33 other (Shalev-Shwartz and Ben-David, 2014; Thirion et al., 2014), investigators should evaluate cluster  
 34 findings and choose the cluster number by relying on a set of complementary cluster validity criteria, such  
 35 as reproducibility and goodness of fit or bias and variance.

Evidently, the discovered set of connectivity-derived clusters only represent hints to candidate brain modules. Their "existence" in neurobiology requires further scrutiny (Eickhoff et al., 2015; Thirion et al., 2014). Nevertheless, such clustering solutions are an important means to narrow down high-dimensional neuroimaging data. Preliminary clustering results broaden the space of research hypotheses that the investigator can articulate. For instance, unexpected discovery of a candidate brain region (cf. Mars et al., 2012; zu Eulenburg et al., 2012) can provide an argument for future experimental investigations. Brain parcellation can thus be viewed as an exploratory unsupervised method outlining relevant structure in neuroimaging data that can subsequently be tested as research hypotheses in targeted future neuroimaging studies on classical inference or out-of-sample generalization.

## Conclusion

A novel scientific fact about the brain is only valid in the context of the complexity restrictions that have been imposed on the studied phenomenon during the investigation (Box, 1976). Tools of the imaging neuroscientist's statistical arsenal can be placed on a continuum between *classical inference* by hypothesis falsification and increasingly used *out-of-sample generalization* by extrapolating complex patterns to independent data (Efron and Hastie, 2016). While null-hypothesis testing has been dominating academic milieus in the empirical sciences and statistics departments for several decades, statistical learning methods are perhaps still more prevalent in data-intensive industries (Breiman, 2001; Henke et al., 2016; Vanderplas, 2013). This sociological segregation may contribute to the existing confusion about the mutual relationship between the ClSt and StLe camps in application domains such as imaging neuroscience. Despite the incongruent historical trajectories and theoretical foundations, both statistical cultures aim at inferential conclusions by extracting new knowledge from data using mathematical models (Friston et al., 2008; Jordan et al., 2013). However, an observed effect in the brain with a statistically significant p value does not in all cases generalize to future brain recordings (Arbabshirani et al., 2017; Shmueli, 2010; Yarkoni and Westfall, 2016). Conversely, a neurobiological effect that can be successfully captured by a learning algorithm as evidenced by out-of-sample generalization does not invariably entail a significant p value when submitted to null-hypothesis testing. The distributional properties of brain data important for high statistical significance and for high prediction accuracy are not identical (Arbabshirani et al., 2017; Efron, 2012; Lo et al., 2015). The goal and permissible conclusions of a neuroscientific investigation are therefore conditioned by the adopted statistical framework (cf. Feyerabend, 1975). Awareness of the *prediction-inference distinction* will be critical to keep pace with the increasing information detail of neuroimaging data repositories (Bzdok and Yeo, 2017; Eickhoff et al., 2016). Ultimately, statistical inference is not a uniquely defined concept.

## 1 Acknowledgments

2 The present paper did not result from isolated contemplations by a single person. Rather, it emerged from  
3 exposure to several thought milieus with different thought styles and opinion systems.

## 5 Funding

6 This work was supported by the Deutsche Forschungsgemeinschaft (DFG, BZ2/2-1, BZ2/3-1, and BZ2/4-1;  
7 International Research Training Group IRTG2150), Amazon AWS Research Grant (2016 and 2017), the  
8 German National Academic Foundation, and the START-Program of the Faculty of Medicine, RWTH Aachen.

## 10 Figures

### 11 Figure 1: Application areas of two statistical paradigms

12 Lists examples of research domains which apply relatively more classical statistics (*blue*) or learning  
13 algorithms (*red*). The co-occurrence of increased computational resources, growing data repositories, and  
14 improving pattern-learning techniques have initiated a shift towards less hypothesis-driven and more  
15 computer-based methodologies. As a broad intuition, researchers in the empirical sciences on the left tend  
16 to use statistics to evaluate a pre-assumed model on the data. Researchers in the application domains on  
17 the right tend to derive a model directly from the data: A new function with potentially many parameters is  
18 created that can predict the output from the input alone without explicit programming model. One of the  
19 key differences becomes apparent when thinking of the neurobiological phenomenon under study as a  
20 black box (Breiman, 2001). ClSt typically aims at modeling the black box by making a set of formal  
21 assumptions about its content, such as the nature of the signal distribution. Gaussian distributional  
22 assumptions have been very useful in many instances to enhance mathematical convenience and, hence,  
23 computational tractability. Instead, StLe takes a brute-force approach to model the output of the black box  
24 (e.g., tell healthy and schizophrenic people apart) from its input (e.g., volumetric brain measurements)  
25 while making a possible minimum of assumptions (Abu-Mostafa et al., 2012). In ClSt the stochastic  
26 processes that generated the data is therefore treated as partly known, whereas in StLe the phenomenon is  
27 treated as complex, largely unknown, and partly unknowable.

### 28 Figure 2: Developments in the history of classical statistics and statistical learning

29 Examples of important inventions in statistical methodology. Roughly, a number of statistical methods  
30 taught in today's textbooks in psychology and medicine have emerged in the first half of the 20th century  
31 (*blue*). Instead, many algorithmic techniques and procedures have emerged in the second half of the 20th  
32 century (*red*). "The postwar era witnessed a massive expansion of statistical methodology, responding to  
33 the data-driven demands fo modern scientific technology." (Efron and Hastie, 2016)

### 34 Figures 3: Key differences in the modeling philosophy of classical statistics and statistical learning

35 Ten modeling intuitions that tend to be relatively more characteristic for classical statistical methods (*blue*)  
36 or pattern-learning methods (*red*). In comparison to ClSt, StLe "is essentially a form of applied statistics  
37 with increased emphasis on the use of computers to statistically estimate complicated functions and a  
38 decreased emphasis on proving confidence intervals around these functions" (Goodfellow et al., 2016).  
39 Broadly, ClSt tends to be more analytical by imposing mathematical rigor on the phenomenon, whereas  
40 StLe tends to be more heuristic by finding useful approximations. In practice, ClSt is probably more often  
41 applied to experimental data, where a set of target variables are systematically controlled by the  
42 investigator and the brain system under studied has been subject to experimental perturbation. Instead,  
43 StLe is probably more often applied to observational data without such structured influence and where the  
44 studied system has been left unperturbed. ClSt fully species the statistical model at the beginning of the  
45 investigation, whereas in StLe there is a bigger emphasis on models that can flexibly adapt to the data (e.g.,  
46 learning algorithms creating decision trees).

#### Figure 4: Key concepts in classical statistics and statistical learning

Schematic with statistical notions that are relatively more associated with classical statistical methods (*left column*) or pattern-learning methods (*right column*). As there is a smooth transition between the classical statistical toolkit and learning algorithms, some notions may be closely associated with both statistical cultures (*middle column*).

#### Figure 5: Key differences between measuring outcomes in classical statistics and statistical learning

Ten intuitions on quantifying statistical modeling outcomes that tend to be relatively more true for classical statistical methods (*blue*) or pattern-learning methods (*red*). ClSt typically yields point estimates and interval estimates (e.g., p values, variances, confidence intervals), whereas StLe frequently outputs a function or a program that can yield point and interval estimates on new observations (e.g., the k-means centroids or a trained classifier's decision function can be applied to new data). In many cases, classical inference is a judgment about an entire data sample, whereas a trained predictive model can obtain quantitative answers from a single data point.

#### References

- Abu-Mostafa, Y.S., Magdon-Ismael, M., Lin, H.T., 2012. Learning from data. AMLBook, California.
- Altman, D.G., Bland, J.M., 1994. Statistics Notes: Diagnostic tests 2: predictive values. *Bmj* 309, 102.
- Amunts, K., Lepage, C., Borgeat, L., Mohlberg, H., Dickscheid, T., Rousseau, M.E., Bludau, S., Bazin, P.L., Lewis, L.B., Oros-Peusquens, A.M., Shah, N.J., Lippert, T., Zilles, K., Evans, A.C., 2013. BigBrain: an ultrahigh-resolution 3D human brain model. *Science* 340, 1472-1475.
- Anderson, D.R., Burnham, K.P., Thompson, W.L., 2000. Null hypothesis testing: problems, prevalence, and an alternative. *The journal of wildlife management*, 912-923.
- Anderson, M.L., 2010. Neural reuse: a fundamental organizational principle of the brain. *Behav Brain Sci* 33, 245-266; discussion 266-313.
- Arbabshirani, M.R., Plis, S., Sui, J., Calhoun, V.D., 2017. Single subject prediction of brain disorders in neuroimaging: Promises and pitfalls. *Neuroimage* 145, 137-165.
- Averbeck, B.B., Latham, P.E., Pouget, A., 2006. Neural correlations, population coding and computation. *Nat Rev Neurosci* 7, 358-366.
- Bach, F., 2014. Breaking the curse of dimensionality with convex neural networks. *arXiv preprint arXiv:1412.8690*.
- Behrens, T.E., Johansen-Berg, H., Woolrich, M.W., Smith, S.M., Wheeler-Kingshott, C.A., Boulby, P.A., Barker, G.J., Sillery, E.L., Sheehan, K., Ciccarelli, O., Thompson, A.J., Brady, J.M., Matthews, P.M., 2003. Non-invasive mapping of connections between human thalamus and cortex using diffusion imaging. *Nat Neurosci* 6, 750-757.
- Bellec, P., Rosa-Neto, P., Lyttelton, O.C., Benali, H., Evans, A.C., 2010. Multi-level bootstrap analysis of stable clusters in resting-state fMRI. *Neuroimage* 51, 1126-1139.
- Bellman, R.E., 1961. Adaptive control processes: a guided tour. Princeton University Press.
- Bengio, Y., 2014. Evolving culture versus local minima. *Growing Adaptive Machines*. Springer, pp. 109-138.
- Bengio, Y., Courville, A., Vincent, P., 2013. Representation learning: A review and new perspectives. *PAMI, IEEE* 35, 1798-1828.

1 Berk, R., Brown, L., Buja, A., Zhang, K., Zhao, L., 2013. Valid post-selection inference. *The Annals of Statistics*  
2 41, 802-837.

3 Berkson, J., 1938. Some difficulties of interpretation encountered in the application of the chi-square test.  
4 *Journal of the American Statistical Association* 33, 526-536.

5 Bishop, C.M., 2006. *Pattern Recognition and Machine Learning*. Springer, Heidelberg.

6 Bishop, C.M., Lasserre, J., 2007. Generative or Discriminative? Getting the Best of Both Worlds. *Bayesian*  
7 *statistics* 8, 3-24.

8 Blei, D.M., Smyth, P., 2017. Science and data science. *Proceedings of the National Academy of Sciences*,  
9 201702076.

10 Box, G.E.P., 1976. Science and statistics. *Journal of the American Statistical Association* 71, 791-799.

11 Breiman, L., 2001. Statistical Modeling: The Two Cultures. *Statistical Science* 16, 199-231.

12 Brodersen, K.H., 2009. Decoding mental activity from neuroimaging data — the science behind mind-  
13 reading. *The New Collection, Oxford* 4, 50-61.

14 Brodersen, K.H., Daunizeau, J., Mathys, C., Chumbley, J.R., Buhmann, J.M., Stephan, K.E., 2013. Variational  
15 Bayesian mixed-effects inference for classification studies. *Neuroimage* 76, 345-361.

16 Brodersen, K.H., Haiss, F., Ong, C.S., Jung, F., Tittgemeyer, M., Buhmann, J.M., Weber, B., Stephan, K.E.,  
17 2011a. Model-based feature construction for multivariate decoding. *Neuroimage* 56, 601-615.

18 Brodersen, K.H., Schofield, T.M., Leff, A.P., Ong, C.S., Lomakina, E.I., Buhmann, J.M., Stephan, K.E., 2011b.  
19 Generative embedding for model-based classification of fMRI data. *PLoS Comput Biol* 7, e1002079.

20 Bühlmann, P., Van De Geer, S., 2011. *Statistics for high-dimensional data: methods, theory and*  
21 *applications*. Springer Science & Business Media.

22 Burnham, K.P., Anderson, D.R., 2014. P values are only an index to evidence: 20th-vs. 21st-century  
23 statistical science. *Ecology* 95, 627-630.

24 Bzdok, D., Eickenberg, M., Grisel, O., Thirion, B., Varoquaux, G., 2015. Semi-Supervised Factored Logistic  
25 Regression for High-Dimensional Neuroimaging Data. *NIPS*, pp. 3330-3338.

26 Bzdok, D., Eickenberg, M., Varoquaux, G., Thirion, B., 2017. Hierarchical Region-Network Sparsity for High-  
27 Dimensional Inference in Brain Imaging. *Information Processing in Medical Imaging (IPMI)*.

28 Bzdok, D., Varoquaux, G., Grisel, O., Eickenberg, M., Poupon, C., Thirion, B., 2016. Formal models of the  
29 network co-occurrence underlying mental operations. *PLoS Comput Biol*, DOI:  
30 10.1371/journal.pcbi.1004994.

31 Bzdok, D., Yeo, B.T.T., 2017. Inference in the age of big data: Future perspectives on neuroscience.  
32 *Neuroimage*.

33 Casella, G., Berger, R.L., 2002. *Statistical inference*. Duxbury Pacific Grove, CA.

34 Chamberlin, T.C., 1890. The Method of Multiple Working Hypotheses. *Science* 15, 92-96.

35 Chambers, J.M., 1993. Greater or lesser statistics: a choice for future research. *Statistics and Computing* 3,  
36 182-184.



1 Choi, Y., Taylor, J., Tibshirani, R., 2014. Selecting the number of principal components: estimation of the  
2 true rank of a noisy matrix. arXiv preprint arXiv:1410.8260.

3 Chow, S.L., 1998. Precis of statistical significance: rationale, validity, and utility. Behav Brain Sci 21, 169-194;  
4 discussion 194-239.

5 Christoff, K., Irving, Z.C., Fox, K.C.R., Spreng, R.N., Andrews-Hanna, J.R., 2016. Mind-wandering as  
6 spontaneous thought: a dynamic framework. Nature Reviews Neuroscience.

7 Chumbley, J.R., Friston, K.J., 2009. False discovery rate revisited: FDR and topological inference using  
8 Gaussian random fields. Neuroimage 44, 62-70.

9 Cleveland, W.S., 2001. Data science: an action plan for expanding the technical areas of the field of  
10 statistics. International statistical review 69, 21-26.

11 Cohen, J., 1977. Statistical power analysis for the behavioral sciences (rev. Lawrence Erlbaum Associates,  
12 Inc.

13 Cohen, J., 1990. Things I have learned (so far). American Psychologist 45, 1304.

14 Cohen, J., 1992. A power primer. Psychological Bulletin 112, 155.

15 Cohen, J., 1994. The Earth Is Round ( $p < .05$ ). American Psychologist 49, 997-1003.

16 Cowles, M., Davis, C., 1982. On the Origins of the .05 Level of Statistical Significance. American Psychologist  
17 37, 553-558.

18 Cox, D.D., Dean, T., 2014. Neural networks and neuroscience-inspired computer vision. Curr Biol 24, R921-  
19 929.

20 Cox, D.R., 1975. A note on data-splitting for the evaluation of significance levels. Biometrika 62, 441-444.

21 Cumming, G., 2009. Inference by eye: reading the overlap of independent confidence intervals. Stat Med  
22 28, 205-220.

23 Davatzikos, C., 2004. Why voxel-based morphometric analysis should be used with great caution when  
24 characterizing group differences. Neuroimage 23, 17-20.

25 Davis, J., Goadrich, M., 2006. The relationship between Precision-Recall and ROC curves. Proceedings of the  
26 23rd international conference on Machine learning. ACM, pp. 233-240.

27 de Brebisson, A., Montana, G., 2015. Deep Neural Networks for Anatomical Brain Segmentation. arXiv  
28 preprint arXiv:1502.02445.

29 de-Wit, L., Alexander, D., Ekroll, V., Wagemans, J., 2016. Is neuroimaging measuring information in the  
30 brain? Psychon Bull Rev, 1-14.

31 Demšar, J., 2006. Statistical comparisons of classifiers over multiple data sets. The Journal of Machine  
32 Learning Research 7, 1-30.

33 Derrfuss, J., Mar, R.A., 2009. Lost in localization: The need for a universal coordinate database. Neuroimage  
34 48, 1-7.

35 Domingos, P., 2012. A Few Useful Things to Know about Machine Learning. Communications of the ACM  
36 55, 78-87.

37 Donoho, D., 2015. 50 years of Data Science. Tukey Centennial workshop.

- 1 Efron, B., 1979. Bootstrap methods: another look at the jackknife. *The Annals of Statistics*, 1-26.
- 2 Efron, B., 2012. *Large-scale inference: empirical Bayes methods for estimation, testing, and prediction*.  
3 Cambridge University Press.
- 4 Efron, B., Hastie, T., 2016. *Computer-Age Statistical Inference*. Cambridge University Press.
- 5 Efron, B., Tibshirani, R.J., 1991. Statistical data analysis in the computer age. *Science* 253, 390-395.
- 6 Efron, B., Tibshirani, R.J., 1994. *An introduction to the bootstrap*. CRC press.
- 7 Eickhoff, S., Turner, J.A., Nichols, T.E., Van Horn, J.D., 2016. Sharing the wealth: Neuroimaging data  
8 repositories. *Neuroimage* 124, 1065–1068.
- 9 Eickhoff, S.B., Bzdok, D., Laird, A.R., Roski, C., Caspers, S., Zilles, K., Fox, P.T., 2011. Co-activation patterns  
10 distinguish cortical modules, their connectivity and functional differentiation. *Neuroimage* 57, 938-949.
- 11 Eickhoff, S.B., Thirion, B., Varoquaux, G., Bzdok, D., 2015. Connectivity-based parcellation: Critique and  
12 implications. *Hum Brain Mapp*.
- 13 Estes, W.K., 1997. On the communication of information by displays of standard errors and confidence  
14 intervals. *Psychon Bull Rev* 4, 330-341.
- 15 Everitt, B.S., 1979. Unresolved Problems in Cluster Analysis. *Biometrics* 35, 169-181.
- 16 Ferguson, C.J., 2009. An effect size primer: A guide for clinicians and researchers. *Professional Psychology:  
17 Research and Practice* 40, 532.
- 18 Feyerabend, P., 1975. *Against Method: Outline of an Anarchist Theory of Knowledge*. New Left  
19 Books, London.
- 20 Fisher, R.A., 1925. *Statistical methods of research workers*. Oliver and Boyd, London.
- 21 Fisher, R.A., 1935. *The design of experiments*. 1935. Oliver and Boyd, Edinburgh.
- 22 Fisher, R.A., Mackenzie, W.A., 1923. Studies in crop variation. II. The manurial response of different potato  
23 varieties. *The Journal of Agricultural Science* 13, 311-320.
- 24 Fithian, W., Sun, D., Taylor, J., 2014. Optimal inference after model selection. *arXiv preprint*  
25 *arXiv:1410.2597*.
- 26 Fleck, L., Schäfer, L., Schnelle, T., 1935. *Entstehung und Entwicklung einer wissenschaftlichen Tatsache*.  
27 Schwabe Basel.
- 28 Fox, P.T., Lancaster, J.L., Laird, A.R., Eickhoff, S.B., 2014. Meta-analysis in human neuroimaging:  
29 computational modeling of large-scale databases. *Annu Rev Neurosci* 37, 409-434.
- 30 Frackowiak, R., Markram, H., 2015. The future of human cerebral cartography: a novel approach. *Philos  
31 Trans R Soc Lond B Biol Sci* 370.
- 32 Freedman, D.A., 1983. A note on screening regression equations. *The American Statistician* 37, 152-155.
- 33 Friedman, J.H., 1998. Data Mining and Statistics: What's the connection? *Computing Science and Statistics*  
34 29, 3-9.
- 35 Friedman, J.H., 2001. The role of statistics in the data revolution? *International statistical review/revue  
36 internationale de Statistique*, 5-10.

- 1 Friman, O., Cedefamn, J., Lundberg, P., Borga, M., Knutsson, H., 2001. Detection of neural activity in  
2 functional MRI using canonical correlation analysis. *Magnetic resonance in medicine* 45, 323-330.
- 3 Friston, K.J., 2006. Statistical parametric mapping: The analysis of functional brain images. Academic Press,  
4 Amsterdam.
- 5 Friston, K.J., 2009. Modalities, modes, and models in functional neuroimaging. *Science* 326, 399-403.
- 6 Friston, K.J., 2012. Ten ironic rules for non-statistical reviewers. *Neuroimage* 61, 1300-1310.
- 7 Friston, K.J., Chu, C., Mourao-Miranda, J., Hulme, O., Rees, G., Penny, W., Ashburner, J., 2008. Bayesian  
8 decoding of brain images. *Neuroimage* 39, 181-205.
- 9 Friston, K.J., Holmes, A.P., Worsley, K.J., Poline, J.P., Frith, C.D., Frackowiak, R.S., 1994. Statistical  
10 parametric maps in functional imaging: a general linear approach. *Hum Brain Mapp* 2, 189-210.
- 11 Friston, K.J., Liddle, P.F., Frith, C.D., Hirsch, S.R., Frackowiak, R.S.J., 1992. The left medial temporal region  
12 and schizophrenia. *Brain* 115, 367-382.
- 13 Friston, K.J., Price, C.J., Fletcher, P., Moore, C., Frackowiak, R.S.J., Dolan, R.J., 1996. The trouble with  
14 cognitive subtraction. *Neuroimage* 4, 97-104.
- 15 Gabrieli, J.D., Ghosh, S.S., Whitfield-Gabrieli, S., 2015. Prediction as a humanitarian and pragmatic  
16 contribution from human cognitive neuroscience. *Neuron* 85, 11-26.
- 17 Genovese, C.R., Lazar, N.A., Nichols, T., 2002. Thresholding of statistical maps in functional neuroimaging  
18 using the false discovery rate. *Neuroimage* 15, 870-878.
- 19 Ghahramani, Z., 2004. Unsupervised learning. *Advanced lectures on machine learning*. Springer, pp. 72-112.
- 20 Ghahramani, Z., 2015. Probabilistic machine learning and artificial intelligence. *Nature* 521, 452-459.
- 21 Gigerenzer, G., 1993. The superego, the ego, and the id in statistical reasoning. *A handbook for data  
22 analysis in the behavioral sciences: Methodological issues*, 311-339.
- 23 Gigerenzer, G., 2004. Mindless statistics. *The Journal of Socio-Economics* 33, 587-606.
- 24 Gigerenzer, G., Murray, D.J., 1987. *Cognition as intuitive statistics*. Hillsdale, NJ: Erlbaum.
- 25 Giraud, C., 2014. *Introduction to high-dimensional statistics*. CRC Press.
- 26 Gläscher, J., Adolphs, R., Damasio, H., Bechara, A., Rudrauf, D., Calamia, M., Paul, L.K., Tranel, D., 2012.  
27 Lesion mapping of cognitive control and value-based decision making in the prefrontal cortex. *Proc Natl  
28 Acad Sci U S A* 109, 14681-14686.
- 29 Golland, P., Fischl, B., 2003. Permutation tests for classification: towards statistical significance in image-  
30 based studies. *Information processing in medical imaging*. Springer, pp. 330-341.
- 31 Goodfellow, I.J., Bengio, Y., Courville, A., 2016. *Deep learning*. MIT Press, USA.
- 32 Goodman, S.N., 1999. Toward evidence-based medical statistics. 1: The P value fallacy. *Annals of internal  
33 medicine* 130, 995-1004.
- 34 Grady, C.L., Haxby, J.V., Schapiro, M.B., Gonzalez-Aviles, A., Kumar, A., Ball, M.J., Heston, L., Rapoport, S.I.,  
35 1990. Subgroups in dementia of the Alzheimer type identified using positron emission tomography. *J  
36 Neuropsychiatry Clin Neurosci* 2, 373-384.

- 1 Greenwald, A.G., 2012. There is nothing so theoretical as a good method. *Perspectives on Psychological*  
2 *Science* 7, 99-108.
- 3 Güçlü, U., van Gerven, M.A.J., 2015. Deep neural networks reveal a gradient in the complexity of neural  
4 representations across the ventral stream. *The Journal of Neuroscience* 35, 10005-10014.
- 5 Guyon, I., Elisseeff, A., 2003. An Introduction to Variable and Feature Selection. *Journal of machine*  
6 *Learning research* 3, 1157-1182.
- 7 Guyon, I., Weston, J., Barnhill, S., Vapnik, V., 2002. Gene selection for cancer classification using support  
8 vector machines. *Machine Learning* 46, 389-422.
- 9 Halkidi, M., Batistakis, Y., Vazirgiannis, M., 2001. On clustering validation techniques. *Journal of Intelligent*  
10 *Information Systems* 17, 107-145.
- 11 Hall, E.T., 1976. *Beyond culture*. Anchor Books ed.
- 12 Handl, J., Knowles, J., Kell, D.B., 2005. Computational cluster validation in post-genomic data analysis.  
13 *Bioinformatics* 21, 3201-3212.
- 14 Hanke, M., Halchenko, Y.O., Oosterhof, N.N., 2015. PyMVPA Manuel. <http://www.pymvpa.org/>.
- 15 Hanson, S.J., Halchenko, Y.O., 2008. Brain Reading Using Full Brain Support VectorMachines for Object  
16 Recognition: There Is No “Face” Identification Area. *Neural Comput* 20, 486-503.
- 17 Hastie, T., Tibshirani, R., Friedman, J., 2001. *The Elements of Statistical Learning*. Springer Series in  
18 *Statistics*, Heidelberg, Germany.
- 19 Hastie, T., Tibshirani, R., Wainwright, M., 2015. *Statistical Learning with Sparsity: The Lasso and*  
20 *Generalizations*. CRC Press.
- 21 Haxby, J.V., Gobbini, M.I., Furey, M.L., Ishai, A., Schouten, J.L., Pietrini, P., 2001. Distributed and overlapping  
22 representations of faces and objects in ventral temporal cortex. *Science* 293, 2425-2430.
- 23 Haynes, J.-D., 2015. A primer on pattern-based approaches to fMRI: Principles, pitfalls, and perspectives.  
24 *Neuron* 87, 257-270.
- 25 Haynes, J.D., Rees, G., 2005. Predicting the orientation of invisible stimuli from activity in human primary  
26 visual cortex. *Nat Neurosci* 8, 686-691.
- 27 Haynes, J.D., Rees, G., 2006. Decoding mental states from brain activity in humans. *Nat Rev Neurosci* 7,  
28 523-534.
- 29 Henke, N., Bughin, J., Chui, M., Manyika, J., Saleh, T., Wiseman, B., Sethupathy, G., 2016. The age of  
30 analytics: Competing in a data-driven world. Technical report, McKinsey Global Institute.
- 31 Hinton, G.E., Salakhutdinov, R.R., 2006. Reducing the dimensionality of data with neural networks. *Science*  
32 313, 504-507.
- 33 Ioannidis, J.P., 2005. Why most published research findings are false. *PLoS med* 2, e124.
- 34 Jain, A.K., Murty, M.N., Flynn, P.J., 1999. Data clustering: a review. *ACN Computing Surveys* 31, 264-323.
- 35 Jamalabadi, H., Alizadeh, S., Schönauer, M., Leibold, C., Gais, S., 2016. Classification based hypothesis  
36 testing in neuroscience: Below-chance level classification rates and overlooked statistical properties of  
37 linear parametric classifiers. *Hum Brain Mapp* 37, 1842-1855.

- 1 James, G., Witten, D., Hastie, T., Tibshirani, R., 2013. An introduction to statistical learning. Springer.
- 2 Jenatton, R., Audibert, J.-Y., Bach, F., 2011. Structured variable selection with sparsity-inducing norms. The  
3 Journal of Machine Learning Research 12, 2777-2824.
- 4 Jordan, M.I., Committee on the Analysis of Massive Data, Committee on Applied and Theoretical Statistics,  
5 Board on Mathematical Sciences and Their Applications, Division on Engineering and Physical Sciences,  
6 National Research Council, 2013. Frontiers in Massive Data Analysis. The National Academies Press,  
7 Washington, D.C.
- 8 Jordan, M.I., Mitchell, T.M., 2015. Machine learning: Trends, perspectives, and prospects. Science 349, 255-  
9 260.
- 10 Kamitani, Y., Sawahata, Y., 2010. Spatial smoothing hurts localization but not information: pitfalls for brain  
11 mappers. Neuroimage 49, 1949-1952.
- 12 Kamitani, Y., Tong, F., 2005. Decoding the visual and subjective contents of the human brain. Nat Neurosci  
13 8, 679-685.
- 14 Kandel, E.R., Markram, H., Matthews, P.M., Yuste, R., Koch, C., 2013. Neuroscience thinks big (and  
15 collaboratively). Nature Reviews Neuroscience 14, 659-664.
- 16 Kelley, K., Preacher, K.J., 2012. On effect size. Psychol Methods 17, 137.
- 17 King, J.R., Dehaene, S., 2014. Characterizing the dynamics of mental representations: the temporal  
18 generalization method. Trends Cogn Sci 18, 203-210.
- 19 Knops, A., Thirion, B., Hubbard, E.M., Michel, V., Dehaene, S., 2009. Recruitment of an area involved in eye  
20 movements during mental arithmetic. Science 324, 1583-1585.
- 21 Kriegeskorte, N., 2011. Pattern-information analysis: from stimulus decoding to computational-model  
22 testing. Neuroimage 56, 411-421.
- 23 Kriegeskorte, N., Goebel, R., Bandettini, P., 2006. Information-based functional brain mapping. Proc Natl  
24 Acad Sci USA 103, 3863-3868.
- 25 Kriegeskorte, N., Lindquist, M.A., Nichols, T.E., Poldrack, R.A., Vul, E., 2010. Everything you never wanted to  
26 know about circular analysis, but were afraid to ask. J Cereb Blood Flow Metab 30, 1551-1557.
- 27 Kriegeskorte, N., Simmons, W.K., Bellgowan, P.S., Baker, C.I., 2009. Circular analysis in systems  
28 neuroscience: the dangers of double dipping. Nat Neurosci 12, 535-540.
- 29 Kurzweil, R., 2005. The singularity is near: When humans transcend biology. Penguin.
- 30 Lake, B.M., Salakhutdinov, R., Tenenbaum, J.B., 2015. Human-level concept learning through probabilistic  
31 program induction. Science 350, 1332-1338.
- 32 LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. Nature 521, 436-444.
- 33 Lemm, S., Blankertz, B., Dickhaus, T., Muller, K.R., 2011. Introduction to machine learning for brain imaging.  
34 Neuroimage 56, 387-399.
- 35 Lieberman, M.D., Berkman, E.T., Wager, T.D., 2009. Correlations in Social Neuroscience Aren't Voodoo:  
36 Commentary on Vul et al. Perspectives on Psychological Science 4.
- 37 Lo, A., Chernoff, H., Zheng, T., Lo, S.H., 2015. Why significant variables aren't automatically good predictors.  
38 Proc Natl Acad Sci U S A 112, 13892-13897.

1 Loftus, J.R., 2015. Selective inference after cross-validation. arXiv preprint arXiv:1511.08866.

2 Logothetis, N.K., Pauls, J., Augath, M., Trinath, T., Oeltermann, A., 2001. Neurophysiological investigation of  
3 the basis of the fMRI signal. *Nature* 412, 150-157.

4 Manyika, J., Chui, M., Brown, B., Bughin, J., Dobbs, R., Roxburgh, C., Byers, A., 2011. Big data: The next  
5 frontier for innovation, competition, and productivity. Technical report, McKinsey Global Institute.

6 Markram, H., 2012. The human brain project. *Sci Am* 306, 50-55.

7 Mars, R.B., Sallet, J., Schuffelgen, U., Jbabdi, S., Toni, I., Rushworth, M.F., 2012. Connectivity-Based  
8 Subdivisions of the Human Right "Temporoparietal Junction Area": Evidence for Different Areas  
9 Participating in Different Cortical Networks. *Cereb Cortex* 22, 1894-1903.

10 Miller, K.L., Alfaro-Almagro, F., Bangerter, N.K., Thomas, D.L., Yacoub, E., Xu, J., Bartsch, A.J., Jbabdi, S.,  
11 Sotiropoulos, S.N., Andersson, J.L.R., 2016. Multimodal population brain imaging in the UK Biobank  
12 prospective epidemiological study. *Nat Neurosci*.

13 Misaki, M., Kim, Y., Bandettini, P.A., Kriegeskorte, N., 2010. Comparison of multivariate classifiers and  
14 response normalizations for pattern-information fMRI. *Neuroimage* 53, 103-118.

15 Moeller, J.R., Strother, S.C., Sidtis, J.J., Rottenberg, D.A., 1987. Scaled subprofile model: a statistical  
16 approach to the analysis of functional patterns in positron emission tomographic data. *J Cereb Blood Flow*  
17 *Metab* 7, 649-658.

18 Mur, M., Bandettini, P.A., Kriegeskorte, N., 2009. Revealing representational content with pattern-  
19 information fMRI--an introductory guide. *Soc Cogn Affect Neurosci* 4, 101-109.

20 Murphy, K.P., 2012. Machine learning: a probabilistic perspective. MIT press.

21 Naselaris, T., Kay, K.N., Nishimoto, S., Gallant, J.L., 2011. Encoding and decoding in fMRI. *Neuroimage* 56,  
22 400-410.

23 Neyman, J., Pearson, E.S., 1933. On the Problem of the most Efficient Tests for Statistical Hypotheses. *Phil.*  
24 *Trans. R. Soc. A* 231, 289-337.

25 Nichols, T.E., 2012. Multiple testing corrections, nonparametric methods, and random field theory.  
26 *Neuroimage* 62, 811-815.

27 Nichols, T.E., Hayasaka, S., 2003. Controlling the familywise error rate in functional neuroimaging: a  
28 comparative review. *Stat Methods Med Res* 12, 419-446.

29 Nichols, T.E., Holmes, A.P., 2002. Nonparametric permutation tests for functional neuroimaging: a primer  
30 with examples. *Hum Brain Mapp* 15, 1-25.

31 Nickerson, R.S., 2000. Null hypothesis significance testing: a review of an old and continuing controversy.  
32 *Psychol Methods* 5, 241-301.

33 Noirhomme, Q., Lesenfans, D., Gomez, F., Soddu, A., Schrouff, J., Garraux, G., Luxen, A., Phillips, C.,  
34 Laureys, S., 2014. Biased binomial assessment of cross-validated estimation of classification accuracies  
35 illustrated in diagnosis predictions. *NeuroImage: Clinical* 4, 687-694.

36 Nuzzo, R., 2014. Scientific method: statistical errors. *Nature* 506, 150-152.

37 Oakes, M., 1986. Statistical Inference: A commentary for the social and behavioral sciences. Wiley, New  
38 York.

1     Passingham, R.E., Stephan, K.E., Kotter, R., 2002. The anatomical basis of functional localization in the  
2     cortex. *Nat Rev Neurosci* 3, 606-616.

3     Pedregosa, F., Eickenberg, M., Ciuciu, P., Thirion, B., Gramfort, A., 2015. Data-driven HRF estimation for  
4     encoding and decoding models. *Neuroimage* 104, 209-220.

5     Pereira, F., Botvinick, M., 2011. Information mapping with pattern classifiers: a comparative study.  
6     *Neuroimage* 56, 476-496.

7     Pereira, F., Mitchell, T., Botvinick, M., 2009. Machine learning classifiers and fMRI: a tutorial overview.  
8     *Neuroimage* 45, 199-209.

9     Pernet, C.R., Chauveau, N., Gaspar, C., Rousselet, G.A., 2011. LIMO EEG: a toolbox for hierarchical Linear  
10    MOdeling of ElectroEncephaloGraphic data. *Comput Intell Neurosci* 2011, 3.

11    Platt, J.R., 1964. Strong Inference: Certain systematic methods of scientific thinking may produce much  
12    more rapid progress than others. *Science* 146, 347-353.

13    Plis, S.M., Hjelm, D.R., Salakhutdinov, R., Allen, E.A., Bockholt, H.J., Long, J.D., Johnson, H.J., Paulsen, J.S.,  
14    Turner, J.A., Calhoun, V.D., 2014. Deep learning for neuroimaging: a validation study. *Front Neurosci* 8.

15    Poldrack, R.A., 2006. Can cognitive processes be inferred from neuroimaging data? *Trends Cogn Sci* 10, 59-  
16    63.

17    Poldrack, R.A., Baker, C.I., Durnez, J., Gorgolewski, K.J., Matthews, P.M., Munafò, M.R., Nichols, T.E., Poline,  
18    J.-B., Vul, E., Yarkoni, T., 2017. Scanning the horizon: towards transparent and reproducible neuroimaging  
19    research. *Nature Reviews Neuroscience*.

20    Poldrack, R.A., Gorgolewski, K.J., 2014. Making big data open: data sharing in neuroimaging. *Nat Neurosci*  
21    17, 1510-1517.

22    Poline, J.-B., Brett, M., 2012. The general linear model and fMRI: does love last forever? *Neuroimage* 62,  
23    871-880.

24    Popper, K., 1935/2005. *Logik der Forschung*, 11th ed. Mohr Siebeck, Tübingen.

25    Powers, D.M., 2011. Evaluation: from precision, recall and F-measure to ROC, informedness, markedness  
26    and correlation.

27    Rosenblatt, F., 1958. The perceptron: a probabilistic model for information storage and organization in the  
28    brain. *Psychol Rev* 65, 386.

29    Rosnow, R.L., Rosenthal, R., 1989. Statistical procedures and the justification of knowledge in psychological  
30    science. *American Psychologist* 44, 1276.

31    Russell, S.J., Norvig, P., 2002. *Artificial intelligence: a modern approach* (International Edition).

32    Samuel, A.L., 1959. Some studies in machine learning using the game of checkers. *IBM Journal of research*  
33    and development 3, 210-229.

34    Saygin, Z.M., Osher, D.E., Koldewyn, K., Reynolds, G., Gabrieli, J.D., Saxe, R.R., 2012. Anatomical  
35    connectivity patterns predict face selectivity in the fusiform gyrus. *Nat Neurosci* 15, 321-327.

36    Scheffé, H., 1959. *The Analysis of Variance*. Wiley, New York.

37    Schmidt, F.L., 1996. Statistical significance testing and cumulative knowledge in psychology: Implications for  
38    training of researchers. *Psychol Methods* 1, 115.

1 Schwartz, Y., Thirion, B., Varoquaux, G., 2013. Mapping paradigm ontologies to and from the brain.  
2 Advances in neural information processing systems, pp. 1673-1681.

3 Shalev-Shwartz, S., Ben-David, S., 2014. Understanding machine learning: From theory to algorithms.  
4 Cambridge University Press.

5 Shmueli, G., 2010. To explain or to predict? Statistical Science, 289-310.

6 Sladek, R., Rocheleau, G., Rung, J., Dina, C., Shen, L., Serre, D., Boutin, P., Vincent, D., Belisle, A., Hadjadj, S.,  
7 2007. A genome-wide association study identifies novel risk loci for type 2 diabetes. Nature 445, 881-885.

8 Smith, S.M., Beckmann, C.F., Andersson, J., Auerbach, E.J., Bijsterbosch, J., Douaud, G., Duff, E., Feinberg,  
9 D.A., Griffanti, L., Harms, M.P., Kelly, M., Laumann, T., Miller, K.L., Moeller, S., Petersen, S., Power, J.,  
10 Salimi-Khorshidi, G., Snyder, A.Z., Vu, A.T., Woolrich, M.W., Xu, J., Yacoub, E., Ugurbil, K., Van Essen, D.C.,  
11 Glasser, M.F., Consortium, W.U.-M.H., 2013. Resting-state fMRI in the Human Connectome Project.  
12 Neuroimage 80, 144-168.

13 Smith, S.M., Matthews, P.M., Jezzard, P., 2001. Functional MRI: an introduction to methods. Oxford  
14 University Press.

15 Smith, S.M., Nichols, T.E., 2009. Threshold-free cluster enhancement: addressing problems of smoothing,  
16 threshold dependence and localisation in cluster inference. Neuroimage 44, 83-98.

17 Stark, C.E., Squire, L.R., 2001. When zero is not zero: the problem of ambiguous baseline conditions in fMRI.  
18 Proc Natl Acad Sci U S A 98, 12760-12766.

19 Taylor, J., Lockhart, R., Tibshirani, R.J., Tibshirani, R., 2014. Exact post-selection inference for forward  
20 stepwise and least angle regression. arXiv preprint arXiv:1401.3889.

21 Taylor, J., Tibshirani, R.J., 2015. Statistical learning and selective inference. Proc Natl Acad Sci U S A 112,  
22 7629-7634.

23 Tenenbaum, J.B., Kemp, C., Griffiths, T.L., Goodman, N.D., 2011. How to grow a mind: Statistics, structure,  
24 and abstraction. Science 331, 1279-1285.

25 Thirion, B., Varoquaux, G., Dohmatob, E., Poline, J.B., 2014. Which fMRI clustering gives good brain  
26 parcellations? Front Neurosci 8, 167.

27 Tibshirani, R., 1996. Regression shrinkage and selection via the lasso. Journal of the Royal Statistical Society.  
28 Series B (Methodological), 267-288.

29 Tukey, J.W., 1962. The future of data analysis. Annals of Statistics 33, 1-67.

30 UK House of Common, S.a.T., 2016. The big data dilemma. Committee on Applied and Theoretical Statistics,  
31 UK.

32 Van Essen, D.C., Ugurbil, K., Auerbach, E., Barch, D., Behrens, T.E., Bucholz, R., Chang, A., Chen, L., Corbetta,  
33 M., Curtiss, S.W., Della Penna, S., Feinberg, D., Glasser, M.F., Harel, N., Heath, A.C., Larson-Prior, L., Marcus,  
34 D., Michalareas, G., Moeller, S., Oostenveld, R., Petersen, S.E., Prior, F., Schlaggar, B.L., Smith, S.M., Snyder,  
35 A.Z., Xu, J., Yacoub, E., Consortium, W.U.-M.H., 2012. The Human Connectome Project: a data acquisition  
36 perspective. Neuroimage 62, 2222-2231.

37 Van Horn, J.D., Toga, A.W., 2014. Human neuroimaging as a "Big Data" science. Brain Imaging Behav 8, 323-  
38 331.

39 Vanderplas, J., 2013. The Big Data Brain Drain: Why Science is in Trouble. Blog "Pythonic Perambulations".



- 1 Vapnik, V.N., 1989. Statistical Learning Theory. Wiley-Interscience, New York.
- 2 Vapnik, V.N., 1996. The nature of statistical learning theory. Springer, New York.
- 3 Vapnik, V.N., Kotz, S., 1982. Estimation of dependences based on empirical data. Springer-Verlag New York.
- 4 Varoquaux, G., Thirion, B., 2014. How machine learning is shaping cognitive neuroimaging. *GigaScience* 3,  
5 28.
- 6 Vogelstein, J.T., Park, Y., Ohyama, T., Kerr, R.A., Truman, J.W., Priebe, C.E., Zlatić, M., 2014. Discovery of  
7 brainwide neural-behavioral maps via multiscale unsupervised structure learning. *Science* 344, 386-392.
- 8 Vul, E., Harris, C., Winkielman, P., Pashler, H., 2008. Voodoo Correlations in Social Neuroscience. *Psychol*  
9 *Sci.*
- 10 Wainwright, M.J., 2014. Structured Regularizers for High-Dimensional Problems: Statistical and  
11 Computational Issues. *Annu. Rev. Stat. Appl* 1, 233-253.
- 12 Wasserman, L., Roeder, K., 2009. High dimensional variable selection. *Annals of Statistics* 37, 2178.
- 13 Wasserstein, R.L., Lazar, N.A., 2016. The ASA's statement on p-values: context, process, and purpose. *Am*  
14 *Stat* 70, 129-133.
- 15 Wolpert, D., 1996. The lack of a priori distinctions between learning algorithms. *Neural Computation* 8,  
16 1341–1390.
- 17 Worsley, K.J., Evans, A.C., Marrett, S., Neelin, P., 1992. A three-dimensional statistical analysis for CBF  
18 activation studies in human brain. *Journal of Cerebral Blood Flow and Metabolism* 12, 900-900.
- 19 Worsley, K.J., Poline, J.-B., Friston, K.J., Evans, A.C., 1997. Characterizing the response of PET and fMRI data  
20 using multivariate linear models. *Neuroimage* 6, 305-319.
- 21 Wu, T.T., Chen, Y.F., Hastie, T., Sobel, E., Lange, K., 2009. Genome-wide association analysis by lasso  
22 penalized logistic regression. *Bioinformatics* 25, 714-721.
- 23 Yamins, D.L., DiCarlo, J.J., 2016. Using goal-driven deep learning models to understand sensory cortex. *Nat*  
24 *Neurosci* 19, 356-365.
- 25 Yarkoni, T., Braver, T.S., 2010. Cognitive neuroscience approaches to individual differences in working  
26 memory and executive control: conceptual and methodological issues. *Handbook of individual differences*  
27 *in cognition*. Springer, pp. 87-107.
- 28 Yarkoni, T., Poldrack, R.A., Nichols, T.E., Van Essen, D.C., Wager, T.D., 2011. Large-scale automated  
29 synthesis of human functional neuroimaging data. *Nat Methods* 8, 665-670.
- 30 Yarkoni, T., Westfall, J., 2016. Choosing prediction over explanation in psychology: Lessons from machine  
31 learning. *Perspectives on Psychological Science*.
- 32 Yeo, B.T., Krienen, F.M., Chee, M.W., Buckner, R.L., 2014. Estimates of segregation and overlap of  
33 functional connectivity networks in the human cerebral cortex. *Neuroimage* 88, 212-227.
- 34 Yeo, B.T., Krienen, F.M., Sepulcre, J., Sabuncu, M.R., Lashkari, D., Hollinshead, M., Roffman, J.L., Smoller,  
35 J.W., Zolke, L., Polimeni, J.R., Fischl, B., Liu, H., Buckner, R.L., 2011. The organization of the human cerebral  
36 cortex estimated by intrinsic functional connectivity. *J Neurophysiol* 106, 1125-1165.
- 37 Yuste, R., 2015. From the neuron doctrine to neural networks. *Nature Reviews Neuroscience* 16, 487-497.

- 1 Zou, H., Hastie, T., 2005. Regularization and variable selection via the elastic net. Journal of the Royal
- 2 Statistical Society: Series B (Statistical Methodology) 67, 301-320.
- 3 zu Eulenburg, P., Caspers, S., Roski, C., Eickhoff, S.B., 2012. Meta-analytical definition and functional
- 4 connectivity of the human vestibular cortex. Neuroimage 60, 162-169.
- 5